# C H A P T E R   2
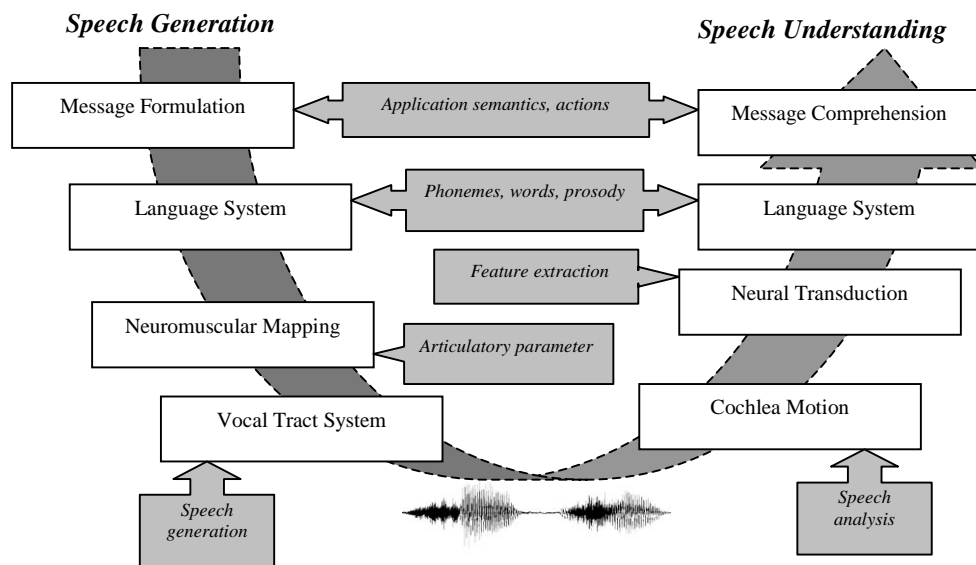
# Spoken Language Structure

$S$poken language is used to communicate information from a speaker to a listener. Speech production and perception are both important components of the speech chain. Speech begins with a thought and intent to communicate in the brain, which activates muscular movements to produce speech sounds. A listener receives it in the auditory system, processing it for conversion to neurological signals the brain can understand. The speaker continuously monitors and controls the vocal organs by receiving his or her own speech as feedback.

Considering the universal components of speech communication as shown in Figure 2.1, the fabric of spoken interaction is woven from many distinct elements. The speech production process starts with the semantic message in a person's mind to be transmitted to the listener via speech. The computer counterpart to the process of message formulation is the application semantics that creates the concept to be expressed. After the message is created, the next step is to convert the message into a sequence of words. Each word consists of a sequence of phonemes that corresponds to the pronunciation of the words. Each sentence also contains a prosodic pattern that denotes the duration of each phoneme, intonation of the sentence, and loudness of the sounds. Once the language system finishes

sentence, and loudness of the sounds. Once the language system finishes the mapping, the talker executes a series of neuromuscular signals. The neuromuscular commands perform articulatory mapping to control the vocal cords, lips, jaw, tongue, and velum, thereby producing the sound sequence as the final output. The speech understanding process works in reverse order. First the signal is passed to the cochlea in the inner ear, which performs frequency analysis as a filter bank. A neural transduction process follows and converts the spectral signal into activity signals on the auditory nerve, corresponding roughly to a feature extraction component. Currently, it is unclear how neural activity is mapped into the language system and how message comprehension is achieved in the brain.



**Figure 2.1** The underlying determinants of speech generation and understanding. The gray boxes indicate the corresponding computer system components for spoken language processing.
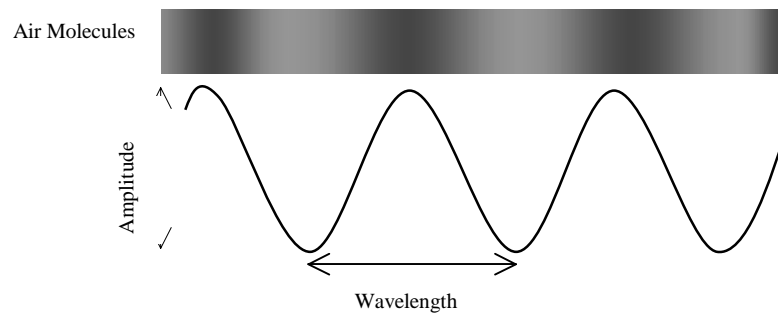
Speech signals are composed of analog sound patterns that serve as the basis for a discrete, symbolic representation of the spoken language – phonemes, syllables, and words. The production and interpretation of these sounds are governed by the syntax and semantics of the language spoken. In this chapter, we take a bottom up approach to introduce the basic concepts from sound to phonetics and phonology. Syllables and words are followed by syntax and semantics, which forms the structure of spoken language processing. The examples in this book are drawn primarily from English, though they are relevant to other languages.

## 2.1. SOUND AND HUMAN SPEECH SYSTEMS

In this Section, we briefly review human speech production and perception systems. We hope spoken language research will enable us to build a computer system that is as good as or better than our own speech production and understanding system.

### 2.1.1. Sound

Sound is a longitudinal pressure wave formed of compressions and rarefactions of air molecules, in a direction parallel to that of the application of energy. Compressions are zones where air molecules have been forced by the application of energy into a tighter-than-usual configuration, and rarefactions are zones where air molecules are less tightly packed. The alternating configurations of compression and rarefaction of air molecules along the path of an energy source are sometimes described by the graph of a sine wave as shown in Figure 2.2. In this representation, crests of the sine curve correspond to moments of maximal compression and troughs to moments of maximal rarefaction.



**Figure 2.2** Application of sound energy causes alternating compression/refraction of air molecules, described by a sine wave. There are two important parameters, amplitude and wavelength, to describe a sine wave. Frequency [cycles/second measured in Hertz (Hz)] is also used to measure of the waveform.

The use of the sine graph in Figure 2.2 is only a notational convenience for charting local pressure variations over time, since sound does not form a transverse wave, and the air particles are just *oscillating in place* along the line of application of energy. The speed of a sound pressure wave in air is approximately $331.5 + 0.6T_c\,m/s$, where $T_c$ is the Celsius temperature.

The amount of work done to generate the energy that sets the air molecules in motion is reflected in the amount of displacement of the molecules from their resting position. This *degree of displacement* is measured as the amplitude of a sound as shown in Figure 2.2. Because of the wide range, it is convenient to measure sound amplitude on a logarithmic scale in *decibels* (dB). A decibel scale is actually a means for comparing two sounds:

$$10\log_{10} \ (P1/P2) \tag{2.1}$$

where $P_1$ and $P_2$ are the two power levels.

Sound pressure level (SPL) is a measure of absolute sound pressure $P$ in dB:

$$SPL(dB) = 20\log_{10}\left(\frac{P}{P_0}\right) \tag{2.2}$$

where the reference 0 dB corresponds to the threshold of hearing, which is $P_0 = 0.0002\mu bar$ for a tone of 1kHz. The speech conversation level at 3 feet is about 60 dB SPL, and a jackhammer's level is about 120 dB SPL. Alternatively, watts/meter$^2$ units are often used to indicate intensity. We can bracket the limits of human hearing as shown in Table 2.1. On the low end, the human ear is quite sensitive. A typical person can detect sound waves having an intensity of $10^{-12}$ W/m$^2$ (the *threshold of hearing* or TOH). This intensity corresponds to a pressure wave affecting a given region by only one-billionth of a centimeter of molecular motion. On the other end, the most intense sound that can be safely detected without suffering physical damage is one billion times more intense than the TOH. 0 dB begins with the TOH and advances logarithmically. The faintest audible sound is arbitrarily assigned a value of 0 dB, and the loudest sounds that the human ear can tolerate are about 120 dB.
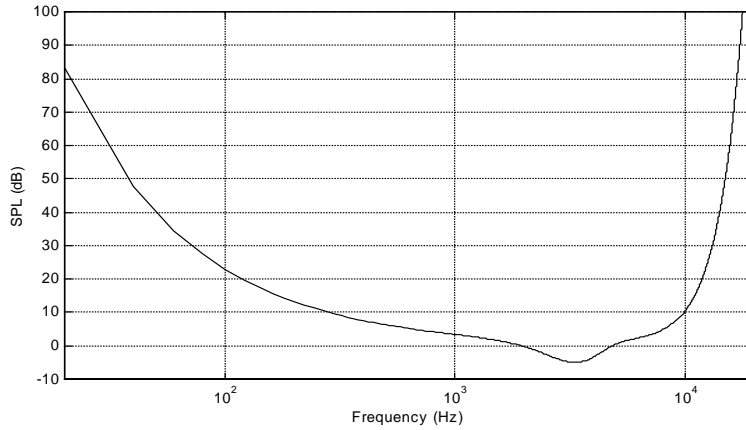
**Table 2.1** Intensity and decibel levels of various sounds.

| Sound | dB Level | Times > TOH |
|---|---|---|
| Threshold of hearing (TOH: $10^{-12}W/m^2$ ) | 0 | $10^0$ |
| Light whisper | 10 | $10^1$ |
| Quiet living room | 20 | $10^2$ |
| Quiet conversation | 40 | $10^4$ |
| Average office | 50 | $10^5$ |
| Normal conversation | 60 | $10^6$ |
| Busy city street | 70 | $10^7$ |
| Acoustic guitar – 1 ft. away | 80 | $10^8$ |
| Heavy truck traffic | 90 | $10^9$ |
| Subway from platform | 100 | $10^{10}$ |
| Power tools | 110 | $10^{11}$ |
| Pain threshold of ear | 120 | $10^{12}$ |
| Airport runway | 130 | $10^{13}$ |
| Sonic boom | 140 | $10^{14}$ |
| Permanent damage to hearing | 150 | $10^{15}$ |
| Jet engine, close up | 160 | $10^{16}$ |
| Rocket engine | 180 | $10^{18}$ |
| Twelve feet. from artillery cannon muzzle ($10^{10}W/m^2$ ) | 220 | $10^{22}$ |

The absolute threshold of hearing is the maximum amount of energy of a pure tone that cannot be detected by a listener in a noise free environment. The absolute threshold of hearing is a function of frequency that can be approximated by

$$T_q(f) = 3.64(f/1000)^{-0.8} - 6.5e^{-0.6(f/1000-3.3)^2} + 10^{-3}(f/1000)^4 \quad (dB \; SPL) \tag{2.3}$$

and is plotted in Figure 2.3.



**Figure 2.3** The sound pressure level (SPL) level in dB of the absolute threshold of hearing as a function of frequency. Sounds below this level are inaudible. Note that below 100 Hz and above 10 kHz this level rises very rapidly. Frequency goes from 20 Hz to 20 kHz and is plotted in a logarithmic scale from Eq. (2.3).

Let's compute how the pressure level varies with distance for a sound wave emitted by a point source located a distance $r$ away. Assuming no energy absorption or reflection, the sound wave of a point source is propagated in a spherical front, such that the energy is the same for the sphere's surface at all radius $r$. Since the surface of a sphere of radius $r$ is $4\pi r^2$, the sound's energy is inversely proportional to $r^2$, so that every time the distance is doubled, the sound pressure level decreases by 6 dB. For the point sound source, the energy ($E$) transported by a wave is proportional to the square of the amplitude ($A$) of the wave and the distance ($r$) between the sound source and the listener:

$$E \propto \frac{A^2}{r^2} \tag{2.4}$$

The typical sound intensity of a speech signal one inch away (close-talking micro-phone) from the talker is 1 Pascal = 10μbar, which corresponds to 94 dB SPL. The typical sound intensity 10 inches away from a talker is 0.1 Pascal = 1μbar, which corresponds to 74dB SPL.
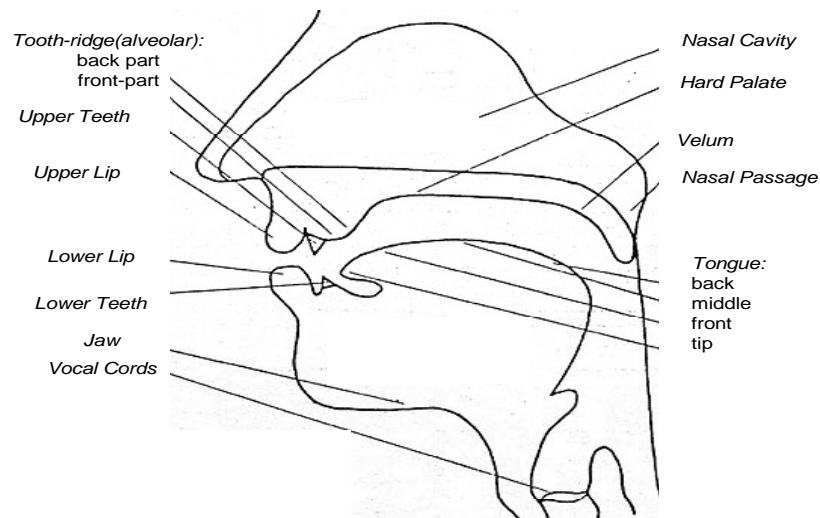
## 2.1.2.    Speech Production

We review here basic human speech production systems, which have influenced research on speech coding, synthesis, and recognition.

### 2.1.2.1.    Articulators

Speech is produced by air-pressure waves emanating from the mouth and the nostrils of a speaker. In most of the world's languages, the inventory of *phonemes,* as discussed in Section 2.2.1, can be split into two basic classes:

> consonants - articulated in the presence of constrictions in the throat or obstructions in the mouth (tongue, teeth, lips) as we speak.

> vowels - articulated without major constrictions and obstructions.

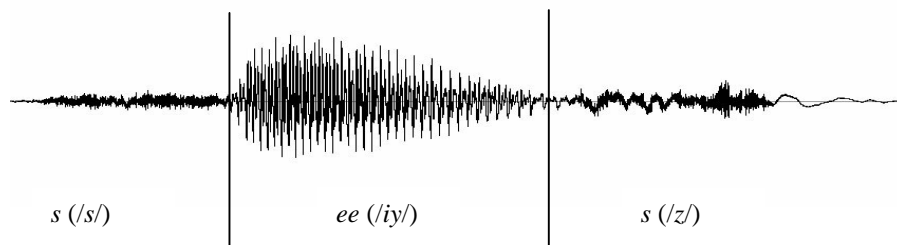**Figure 2.4** A schematic diagram of the human speech production apparatus.

The sounds can be further partitioned into subgroups based on certain articulatory properties. These properties derive from the anatomy of a handful of important articulators and the places where they touch the boundaries of the human vocal tract. Additionally, a large number of muscles contribute to articulator positioning and motion. We restrict ourselves to a schematic view of only the major articulators, as diagrammed in Figure 2.4. The gross components of the speech production apparatus are the lungs, trachea, larynx (organ of voice production), pharyngeal cavity (throat), oral and nasal cavity. The pharyngeal and oral

cavities are typically referred to as the vocal tract, and the nasal cavity as the nasal tract. As illustrated in Figure 2.4, the human speech production apparatus consists of:

*Lungs*: source of air during speech.

*Vocal cords (larynx):* when the vocal folds are held close together and oscillate against one another during a speech sound, the sound is said to be *voiced*. When the folds are too slack or tense to vibrate periodically, the sound is said to be *unvoiced*. The place where the vocal folds come together is called the *glottis*.

*Velum (Soft Palate)*: operates as a *valve*, opening to allow passage of air (and thus resonance) through the nasal cavity. Sounds produced with the flap open include *m* and *n*.

*Hard palate*: a long relatively hard surface at the roof inside the mouth, which, when the tongue is placed against it, enables consonant articulation.

*Tongue*: flexible articulator, shaped away from the palate for vowels, placed close to or on the palate or other hard surfaces for consonant articulation.

*Teeth*: another place of articulation used to brace the tongue for certain consonants.

*Lips*: can be rounded or spread to affect vowel quality, and closed completely to stop the oral air flow in certain consonants (*p, b, m*).
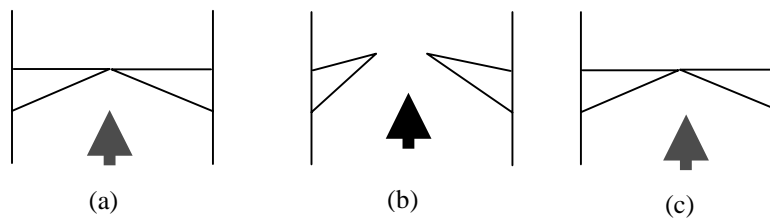
## 2.1.2.2.    The Voicing Mechanism

The most fundamental distinction between sound types in speech is the voiced/voiceless distinction. Voiced sounds, including vowels, have in their time and frequency structure a roughly regular pattern that voiceless sounds, such as consonants like *s*, lack. Voiced sounds typically have more energy as shown in Figure 2.5. We see here the waveform of the word *sees*, which consists of three phonemes: an unvoiced consonant /*s*/, a vowel /*iy*/ and, a voiced consonant /*z*/.

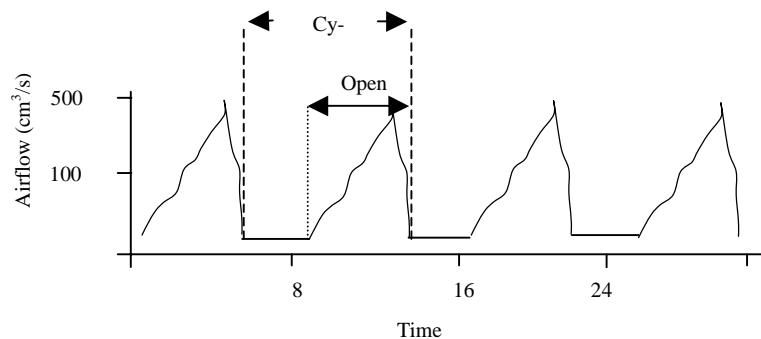$$s\ (/s/) \qquad\qquad ee\ (/iy/) \qquad\qquad s\ (/z/)$$

**Figure 2.5** Waveform of *sees*, showing a voiceless phoneme /*s*, followed by a voiced sound, the vowel /*iy*/. The final sound, /*z*/, is a type of voiced consonant.

What in the speech production mechanism creates this fundamental distinction? When the vocal folds vibrate during phoneme articulation, the phoneme is considered voiced; oth-

erwise it is unvoiced. Vowels are voiced throughout their duration. The distinct vowel *tim-bres* are created by using the tongue and lips to shape the main oral resonance cavity in different ways. The vocal folds vibrate at slower or faster rates, from as low as 60 cycles per second (Hz) for a large man, to as high as 300 Hz or higher for a small woman or child. The rate of cycling (opening and closing) of the vocal folds in the larynx during phonation of voiced sounds is called the *fundamental frequency*. This is because it sets the periodic baseline for all higher-frequency harmonics contributed by the pharyngeal and oral resonance cavities above. The fundamental frequency also contributes more than any other single factor to the perception of *pitch* (the semi-musical rising and falling of voice tones) in speech.



(a)                            (b)                            (c)

**Figure 2.6** Vocal fold cycling at the larynx. (a) Closed with sub-glottal pressure buildup; (b) trans-glottal pressure differential causing folds to blow apart; (c) pressure equalization and tissue elasticity forcing temporary reclosure of vocal folds, ready to begin next cycle.



**Figure 2.7** Waveform showing air flow during laryngeal cycle.

The glottal cycle is illustrated in Figure 2.6. At stage (a), the vocal folds are closed and the air stream from the lungs is indicated by the arrow. At some point, the air pressure on the underside of the barrier formed by the vocal folds increases until it overcomes the resistance of the vocal fold closure and the higher air pressure below blows them apart (b). However, the tissues and muscles of the larynx and the vocal folds have a natural elasticity which tends to make them fall back into place rapidly, once air pressure is temporarily equalized (c). The successive airbursts resulting from this process are the source of energy for all voiced sounds. The time for a single open-close cycle depends on the stiffness and size of
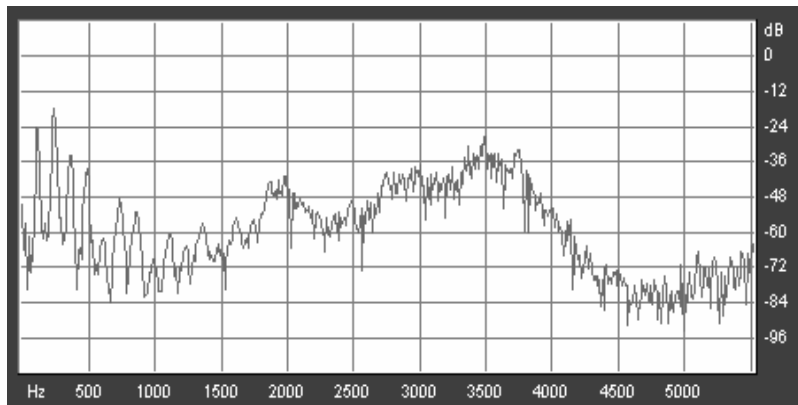
the vocal folds and the amount of subglottal air pressure. These factors can be controlled by a speaker to raise and lower the perceived frequency or pitch of a voiced sound.

The waveform of air pressure variations created by this process can be described as a periodic flow, in cubic centimeters per second (after [15]). As shown in Figure 2.7, during the time bracketed as *one cycle*, there is no air flow during the initial closed portion. Then as the glottis opens (open phase), the volume of air flow becomes greater. After a short peak, the folds begin to resume their original position and the air flow declines until complete closure is attained, beginning the next cycle. A common measure is the number of such cycles per second (Hz), or the fundamental frequency (*F0*). Thus the fundamental frequency for the waveform in Figure 2.7 is about 120 Hz.
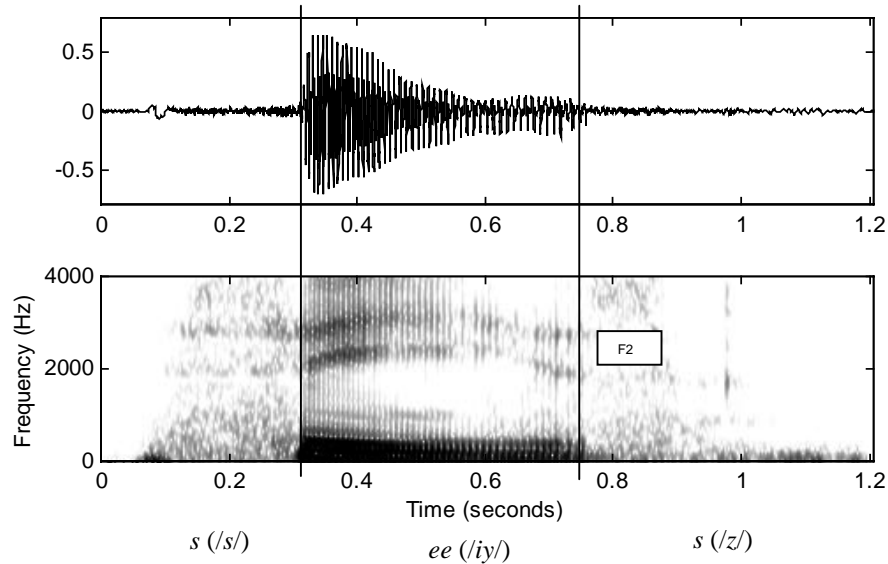
### 2.1.2.3. Spectrograms and Formants

Since the glottal wave is periodic, consisting of fundamental frequency (*F0*) and a number of harmonics (integral multiples of *F0*), it can be analyzed as a sum of sine waves as discussed in Chapter 5. The resonances of the vocal tract (above the glottis) are excited by the glottal energy. Suppose, for simplicity, we regard the vocal tract as a straight tube of uniform cross-sectional area, closed at the glottal end, open at the lips. When the shape of the vocal tract changes, the resonances change also. Harmonics near the resonances are emphasized, and, in speech, the resonances of the cavities that are typical of particular articulator configurations (e.g., the different vowel timbres) are called *formants*. The vowels in an actual speech waveform can be viewed from a number of different perspectives, emphasizing either a *cross-sectional* view of the harmonic responses at a single moment, or a longer-term view of the formant track evolution over time. The actual spectral analysis of a vowel at a single time-point, as shown in Figure 2.8, gives an idea of the uneven distribution of energy in resonances for the vowel /iy/ in the waveform for *see*, which is shown in Figure 2.5.



**Figure 2.8** A spectral analysis of the vowel /iy/, showing characteristically uneven distribution of energy at different frequencies.

Another view of *sees* of Figure 2.5, called a spectrogram, is displayed in the lower part of Figure 2.9. It shows a long-term frequency analysis, comparable to a complete series of single time-point *cross sections* (such as that in Figure 2.8) ranged alongside one another in time and viewed from *above*.



**Figure 2.9** The spectrogram representation of the speech waveform *sees* (approximate phone boundaries are indicated with heavy vertical lines).

In the spectrogram of Figure 2.9, the darkness or lightness of a band indicates the relative amplitude or energy present at a given frequency. The dark horizontal bands show the formants, which are harmonics of the fundamental at natural resonances of the vocal tract cavity position for the vowel */iy/* in *see*. The mathematical methods for deriving analyses and representations such as those illustrated above are covered in Chapters 5 and 6.

## 2.1.3.    Speech Perception

There are two major components in the auditory perception system: the peripheral auditory organs (ears) and the auditory nervous system (brain). The ear processes an acoustic pressure signal by first transforming it into a mechanical vibration pattern on the basilar membrane, and then representing the pattern by a series of pulses to be transmitted by the auditory nerve. Perceptual information is extracted at various stages of the auditory nervous system. In this section we focus mainly on the auditory organs.

## 2.1.3.1.    Physiology of the Ear

The human ear, as shown in Figure 2.10, has three sections: the outer ear, the middle ear and the inner ear. The outer ear consists of the external visible part and the external auditory canal that forms a tube along which sound travels. This tube is about 2.5 *cm* long and is covered by the eardrum at the far end. When air pressure variations reach the eardrum from the outside, it vibrates, and transmits the vibrations to bones adjacent to its opposite side. The vibration of the eardrum is at the same frequency (alternating compression and rarefaction) as the incoming sound pressure wave. The middle ear is an air-filled space or cavity about 1.3 cm across, and about 6 cm$^3$ volume. The air travels to the middle ear cavity along the tube (when opened) that connects the cavity with the nose and throat. The oval window shown in Figure 2.10 is a small membrane at the bony interface to the inner ear (cochlea). Since the cochlear walls are bony, the energy is transferred by mechanical action of the stapes into an impression on the membrane stretching over the oval window.



**Figure 2.10** The structure of the peripheral auditory system with the outer, middle, and inner ear (after Lindsey and Norman [26]).

The relevant structure of the inner ear for sound perception is the cochlea, which communicates directly with the auditory nerve, conducting a representation of sound to the brain. The cochlea is a spiral tube about 3.5 cm long, which coils about 2.6 times. The spiral is divided, primarily by the basilar membrane running lengthwise, into two fluid-filled chambers. The cochlea can be roughly regarded as a filter bank, whose outputs are ordered by location, so that a frequency-to-place transformation is accomplished. The filters closest to the cochlear base respond to the higher frequencies, and those closest to its apex respond to the lower.

### 2.1.3.2.     **Physical vs Perceptual Attributes**

In psychoacoustics, a basic distinction is made between the perceptual attributes of a sound, especially a speech sound, and the measurable physical properties that characterize it. Each of the perceptual attributes, as listed in Table 2.2, seems to have a strong correlation with one main physical property, but the connection is complex, because other physical properties of the sound may affect perception in complex ways.

**Table 2.2** Relation between perceptual and physical attributes of sound.

| Physical Quantity | Perceptual Quality |
|---|---|
| Intensity | Loudness |
| Fundamental frequency | Pitch |
| Spectral shape | Timbre |
| Onset/offset time | Timing |
| Phase difference in binaural hearing | Location |



**Figure 2.11** Equal-loudness curves indicate that the response of the human hearing mechanism is a function of frequency and loudness levels. This relationship again illustrates the difference between physical dimensions and psychological experience (after ISO 226).

Although sounds with a greater intensity level usually sound louder, the sensitivity of the ear varies with the frequency and the quality of the sound. One fundamental divergence between physical and perceptual qualities is the phenomenon of non-uniform *equal loudness* perception of tones of varying frequencies. In general, tones of differing pitch have different

inherent *perceived loudness*. The sensitivity of the ear varies with the frequency and the quality of the sound. The graph of equal loudness contours adopted by ISO is shown in Figure 2.11. These curves demonstrate the relative insensitivity of the ear to sounds of low frequency at moderate to low intensity levels. Hearing sensitivity reaches a maximum around 4000 Hz, which is near the first resonance frequency of the outer ear canal, and peaks again around 13 kHz, the frequency of the second resonance [38].

*Pitch* is indeed most closely related to the fundamental frequency. The higher the fundamental frequency, the higher the pitch we perceive. However, discrimination between two pitches depends on the frequency of the lower pitch. Perceived pitch will change as intensity is increased and frequency is kept constant.

In another example of the non-identity of acoustic and perceptual effects, it has been observed experimentally that when the ear is exposed to two or more different tones, it is a common experience that one tone may *mask* the others. Masking is probably best explained as an upward shift in the hearing threshold of the weaker tone by the louder tone. Pure tones, complex sounds, narrow and broad bands of noise all show differences in their ability to mask other sounds. In general, pure tones close together in frequency mask each other more than tones widely separated in frequency. A pure tone masks tones of higher frequency more effectively than tones of lower frequency. The greater the intensity of the masking tone, the broader the range of the frequencies it can mask [18, 31].

Binaural listening greatly enhances our ability to sense the direction of the sound source. The sense of localization attention is mostly focused on side-to-side discrimination or *lateralization*. Time and intensity cues have different impacts for low frequency and high frequency, respectively. Low-frequency sounds are lateralized mainly on the basis of interaural time difference, whereas high-frequency sounds are localized mainly on the basis of interaural intensity differences [5].

Finally, an interesting perceptual issue is the question of distinctive voice quality. Speech from different people sounds different. Partially this is due to obvious factors, such as differences in characteristic fundamental frequency caused by, for example, the greater mass and length of adult male vocal folds as opposed to female. But there are more subtle effects as well. In psychoacoustics, the concept of *timbre* (of a sound or instrument) is defined as that attribute of auditory sensation by which a subject can judge that two sounds similarly presented and having the same loudness and pitch are dissimilar. In other words, when all the easily measured differences are controlled, the remaining perception of difference is ascribed to timbre. This is heard most easily in music, where the same note in the same octave played for the same duration on a violin sounds different from a flute. The timbre of a sound depends on many physical variables including a sound's spectral power distribution, its temporal envelope, rate and depth of amplitude or frequency modulation, and the degree of inharmonicity of its harmonics.

### 2.1.3.3.     Frequency Analysis

Researchers have undertaken psychoacoustic experimental work to derive frequency scales that attempt to model the natural response of the human perceptual system, since the cochlea of the inner ear acts as a spectrum analyzer. The complex mechanism of the inner ear and

auditory nerve implies that the perceptual attributes of sounds at different frequencies may not be entirely simple or linear in nature. It is well known that the western musical pitch is described in *octaves*[1] and *semi-tones*[2]. The perceived musical pitch of complex tones is basically proportional to the logarithm of frequency. For complex tones, the just noticeable difference for frequency is essentially constant on the octave/semi-tone scale. Musical pitch scales are used in prosodic research (on speech intonation contour generation).

**Table 2.3** The Bark frequency scale.

| Bark Band # | Edge (Hz) | Center (Hz) |
|:---:|:---:|:---:|
| 1 | 100 | 50 |
| 2 | 200 | 150 |
| 3 | 300 | 250 |
| 4 | 400 | 350 |
| 5 | 510 | 450 |
| 6 | 630 | 570 |
| 7 | 770 | 700 |
| 8 | 920 | 840 |
| 9 | 1080 | 1000 |
| 10 | 1270 | 1170 |
| 11 | 1480 | 1370 |
| 12 | 1720 | 1600 |
| 13 | 2000 | 1850 |
| 14 | 2320 | 2150 |
| 15 | 2700 | 2500 |
| 16 | 3150 | 2900 |
| 17 | 3700 | 3400 |
| 18 | 4400 | 4000 |
| 19 | 5300 | 4800 |
| 20 | 6400 | 5800 |
| 21 | 7700 | 7000 |
| 22 | 9500 | 8500 |
| 23 | 12000 | 10500 |
| 24 | 15500 | 13500 |

AT&T Bell Labs has contributed many influential discoveries in hearing, such as critical band and articulation index, since the turn of the 20th century [3]. Fletcher's work [14] pointed to the existence of critical bands in the cochlear response. Critical bands are of great importance in understanding many auditory phenomena such as perception of loudness, pitch, and timbre. The auditory system performs frequency analysis of sounds into their

---

[1] A tone of frequency $f_1$ is said to be an octave above a tone with frequency $f_2$ if and only if $f_1 = 2 f_2$.

[2] There are 12 semitones in one octave, so a tone of frequency $f_1$ is said to be a semitone above a tone with frequency $f_2$ if and only if $f_1 = 2^{1/12} f_2 = 1.05946 f_2$.

component frequencies. The cochlea acts as if it were made up of overlapping filters having bandwidths equal to the critical bandwidth. One class of critical band scales is called *Bark frequency scale*. It is hoped that by treating spectral energy over the Bark scale, a more natural fit with spectral information processing in the ear can be achieved. The Bark scale ranges from 1 to 24 Barks, corresponding to 24 critical bands of hearing as shown in Table 2.3. As shown in Figure 2.12, the perceptual resolution is finer in the lower frequencies. It should be noted that the ear's critical bands are continuous, and a tone of any audible frequency always finds a critical band centered on it. The Bark frequency $b$ can be expressed in terms of the linear frequency (in Hz) by

$$b(f) = 13 \arctan(0.00076f) + 3.5 * \arctan\left((f/7500)^2\right) \quad (Bark) \tag{2.5}$$



**Figure 2.12** The center frequency of 24 Bark frequency filters as illustrated in Table 2.3.

Another such perceptually motivated scale is the mel frequency scale [41], which is linear below 1 kHz, and logarithmic above, with equal numbers of samples taken below and above 1 kHz. The mel scale is based on experiments with simple tones (sinusoids) in which subjects were required to divide given frequency ranges into four perceptually equal intervals or to adjust the frequency of a stimulus tone to be half as high as that of a comparison tone. One mel is defined as one thousandth of the pitch of a 1 kHz tone. As with all such attempts, it is hoped that the mel scale more closely models the sensitivity of the human ear than a purely linear scale and provides for greater discriminatory capability between speech segments. Mel-scale frequency analysis has been widely used in modern speech recognition systems. It can be approximated by:

$$B(f) = 1125 \ln(1 + f/700) \tag{2.6}$$

The mel scale is plotted in Figure 2.13 together with the Bark scale and the bilinear transform (see Chapter 6).

**Figure 2.13** Frequency warping according to the Bark scale, ERB scale, mel-scale and bilinear transform for $\alpha = 0.6$ : linear frequency in the *x*-axis and normalized frequency in the *y*-axis.

A number of techniques in the modern spoken language system, such as cepstral analysis, and dynamic feature, have benefited tremendously from perceptual research as discussed throughout this book.

### 2.1.3.4. Masking

Frequency masking is a phenomenon under which one sound cannot be perceived if another sound close in frequency has a high enough level. The first sound *masks* the other one. Frequency-masking levels have been determined empirically, with complicated models that take into account whether the masker is a tone or noise, the masker's level, and other considerations.

We now describe a phenomenon known as *tone-masking noise*. It has been determined empirically that noise with energy $E_N$ (dB) masks a tone at bark frequency *b* if the tone's energy is below the threshold

$$T_T(b) = E_N - 6.025 - 0.275i + S_m(b) \quad (dB \; SPL) \tag{2.7}$$

where *K* has been typically set between 3 and 5 dB, and where the *spread-of-masking* function $S_m(b)$ is given by

$$S_m(b) = 15.81 + 7.5(b + 0.474) - 17.5\sqrt{1 + (b + 0.474)^2} \quad (dB) \tag{2.8}$$

We now describe a phenomenon known as *noise-masking tone.* It has been determined empirically that a tone at critical band number $i$ with energy $E_T$ (dB) masks noise at bark frequency $b$ if the noise energy is below the threshold

$$T_N(b) = E_T - 2.025 - 0.175i + S_m(b) \quad (dB\ SPL) \tag{2.9}$$

Masking thresholds are commonly referred to in the literature as Bark scale functions of *just noticeable distortion* (JND). Equation (2.8) can be approximated by a triangular spreading function that has slopes of +25 and –10 dB per Bark, as shown in Figure 2.14.



**Figure 2.14** Contribution of frequency bin $i$ to the masked threshold $S_m(b)$.

In Figure 2.15 we show both the threshold of hearing and the masked threshold of a tone at 1 kHz with a 69 dB SPL. The combined masked threshold is the sum of the two in the linear domain

$$T(f) = 10\log_{10}\left(10^{0.1T_h(f)} + 10^{0.1T_T(f)}\right) \tag{2.10}$$

which is approximately the largest of the two.



**Figure 2.15** Absolute Threshold of Hearing and Spread of Masking threshold for a 1 kHz sinewave masker with a 69 dB SPL. The overall masked threshold is approximately the largest of the two thresholds.

In addition to frequency masking, there is a phenomenon called temporal masking by which a sound too close in time to another sound cannot be perceived. Whereas premasking tends to last about 5 ms, postmasking can last from 50 to 300 ms. Temporal masking level of a masker with a uniform level starting at 0 ms and lasting 200 ms is shown in Figure 2.16.



**Figure 2.16** Temporal masking level of a masker with a uniform level starting at 0 ms and lasting 200 ms.

## 2.2.    PHONETICS AND PHONOLOGY

We now discuss basic phonetics and phonology needed for spoken language processing. *Phonetics* refers to the study of speech sounds and their production, classification, and transcription. *Phonology* is the study of the distribution and patterning of speech sounds in a language and of the tacit rules governing pronunciation.

### 2.2.1.    Phonemes

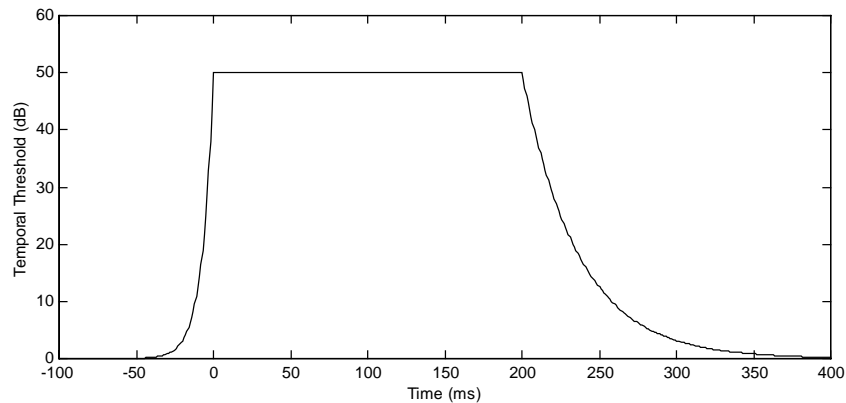Linguist Ferdinand de Saussere (1857-1913) is credited with the observation that the relation between a sign and the object signified by it is arbitrary. The same concept, a certain yellow and black flying social insect, has the sign *honeybee* in English and *mitsubachi* in Japanese. There is no particular relation between the various pronunciations and the meaning, nor do these pronunciations per se describe the bee's characteristics in any detail. For phonetics, this means that the speech sounds described in this chapter have no inherent meaning, and should be randomly distributed across the lexicon, except as affected by extraneous historical or etymological considerations. The sounds are just a set of arbitrary effects made available by human vocal anatomy. You might wonder about this theory when you observe, for example, the number of words beginning with *sn* that have to do with nasal functions in English: *sneeze, snort, sniff, snot, snore, snuffle*, etc. But Saussere's observation is generally true, except for obvious onomatopoetic (sound) words like *buzz*.

Like fingerprints, every speaker's vocal anatomy is unique, and this makes for unique vocalizations of speech sounds. Yet language communication is based on commonality of form at the perceptual level. To allow discussion of the commonalities, researchers have identified certain gross characteristics of speech sounds that are adequate for description and classification of words in dictionaries. They have also adopted various systems of notation to represent the subset of phonetic phenomena that are crucial for meaning.

As an analogy, consider the system of computer coding of text characters. In such systems, the *character* is an abstraction, e.g. the Unicode character U+0041. The identifying property of this character is its Unicode name *LATIN CAPITAL LETTER A*. This is a genuine abstraction; no particular realization is necessarily specified. As the Unicode 2.1 standard [1] states:

*The Unicode Standard does not define glyph images. The standard defines how characters are interpreted, not how glyphs are rendered. The software or hardware-rendering engine of a computer is responsible for the appearance of the characters on the screen. The Unicode Standard does not specify the size, shape, nor orientation of on-screen characters.*

Thus, the U+0041 character can be realized differently for different purposes, and in different sizes with different fonts:

U+0041➜ A, ₐ ₐ A, A, …

The realizations of the character U+0041 are called glyphs, and there is no distinguished uniquely correct glyph for U+0041. In speech science, the term *phoneme* is used to denote any of the minimal units of speech sound in a language that can serve to distinguish one word from another. We conventionally use the term *phone* to denote a phoneme's acoustic realization. In the example given above, U+0041 corresponds to a phoneme and the various fonts correspond to the phone. For example, English phoneme /t/ have two very different acoustic realizations in the words *sat* and *meter*. You had better treat them as two different phones if you want to build a spoken language system. We will use the terms *phone* or *phoneme* interchangeably to refer to the speaker-independent and context-independent units of meaningful sound contrast. Table 2.4 shows a complete list of phonemes used in American English. The set of phonemes will differ in realization across individual speakers. But phonemes will always function systematically to differentiate meaning in words, just as the phoneme /p/ signals the word *pat* as opposed to the similar-sounding but distinct *bat*. The important contrast distinguishing this pair of words is /p/ vs. /b/.

In this section we concentrate on the basic qualities that *define and differentiate abstract phonemes*. In Section 2.2.1.3 below we consider *why and how phonemes vary* in their actual realizations by different speakers and in different contexts.

**Table 2.4** English phonemes used for typical spoken language systems.

| Phonemes | Word Examples | Description |
|---|---|---|
| *ih* | *f**i**ll, h**i**t, l**i**d* | front close unrounded (lax) |
| *ae* | ***a**t, c**a**rry, g**a**s* | front open unrounded (tense) |
| *aa* | *f**a**ther, **ah**, c**a**r* | back open unrounded |
| *ah* | *c**u**t, b**u**d, **u**p* | open-mid back unrounded |
| *ao* | *d**o**g, **law**n, c**au**ght* | open-mid back round |
| *ay* | *t**ie**, **i**ce, b**i**te* | diphthong with quality: aa + ih |
| *ax* | *ag**o**, c**o**mply* | central close mid (schwa) |
| *ey* | ***a**te, d**ay**, t**a**pe* | front close-mid unrounded (tense) |
| *eh* | *p**e**t, b**e**rry, t**e**n* | front open-mid unrounded |
| *er* | *t**ur**n, f**ur**, met**er*** | central open-mid unrounded rhoti- |
| *ow* | *g**o**, **ow**n, t**o**ne* | back close-mid rounded |
| *aw* | *f**ou**l, h**ow**, **ou**r* | diphthong with quality: aa + uh |
| *oy* | *t**oy**, c**oi**n, **oi**l* | diphthong with quality: ao + ih |
| *uh* | *b**oo**k, p**u**ll, g**oo**d* | back close-mid unrounded (lax) |
| *uw* | *t**oo**l, cr**ew**, m**oo*** | back close round |
| *b* | ***b**ig, a**b**le, ta**b*** | voiced bilabial plosive |
| *p* | ***p**ut, o**p**en, ta**p*** | voiceless bilabial plosive |
| *d* | ***d**ig, i**d**ea, wa**d*** | voiced alveolar plosive |
| *t* | ***t**alk, sa**t*** | voiceless alveolar plosive & |
| *t* | *me**t**er* | alveolar flap |
| *g* | ***g**ut, an**g**le, ta**g*** | voiced velar plosive |
| *k* | ***c**ut, **k**en, ta**k**e* | voiceless velar plosive |
| *f* | ***f**ork, a**f**ter, i**f*** | voiceless labiodental fricative |
| *v* | ***v**at, o**v**er, ha**v**e* | voiced labiodental fricative |
| *s* | ***s**it, ca**s**t, to**ss*** | voiceless alveolar fricative |
| *z* | ***z**ap, la**z**y, ha**z**e* | voiced alveolar fricative |
| *th* | ***th**in, no**th**ing, tru**th*** | voiceless dental fricative |
| *dh* | ***Th**en, fa**th**er, scy**the*** | voiced dental fricative |
| *sh* | ***sh**e, cu**sh**ion, wa**sh*** | voiceless postalveolar fricative |
| *zh* | *genre, a**z**ure* | voiced postalveolar fricative |
| *l* | ***l**id* | alveolar lateral approximant |
| *l* | *e**l**bow, sai**l*** | velar lateral approximant |
| *r* | ***r**ed, pa**r**t, fa**r*** | retroflex approximant |
| *y* | ***y**acht,  **y**ard* | palatal sonorant glide |
| *w* | ***w**ith, a**w**ay* | labiovelar sonorant glide |
| *hh* | ***h**elp, a**h**ead, **h**otel* | voiceless glottal fricative |
| *m* | ***m**at, a**m**id, ai**m*** | bilabial nasal |
| *n* | ***n**o, e**n**d, pa**n*** | alveolar nasal |
| *ng* | *si**ng**, a**ng**er* | velar nasal |
| *ch* | ***ch**in, ar**ch**er, mar**ch*** | voiceless alveolar affricate: t + sh |
| *jh* | ***j**oy, a**g**ile, ed**ge*** | voiced alveolar affricate: d + zh |

## 2.2.1.1. Vowels

The tongue shape and positioning in the oral cavity do not form a major constriction of air flow during vowel articulation. However, variations of tongue placement give each vowel its distinct character by changing the resonance, just as different sizes and shapes of bottles give rise to different acoustic effects when struck. The primary energy entering the pharyngeal and oral cavities in vowel production vibrates at the fundamental frequency. The major resonances of the oral and pharyngeal cavities for vowels are called F1 and F2 - the first and second formants, respectively. They are determined by tongue placement and oral tract shape in vowels, and they determine the characteristic timbre or quality of the vowel.

The relationship of F1 and F2 to one another can be used to describe the English vowels. While the shape of the complete vocal tract determines the spectral outcome in a complex, nonlinear fashion, generally F1 corresponds to the back or pharyngeal portion of the cavity, while F2 is determined more by the size and shape of the oral portion, forward of the major tongue extrusion. This makes intuitive sense - the cavity from the glottis to the tongue extrusion is longer than the forward part of the oral cavity, thus we would expect its resonance to be lower. In the vowel of *see*, for example, the tongue extrusion is far forward in the mouth, creating an exceptionally long rear cavity, and correspondingly low F1. The forward part of the oral cavity, at the same time, is extremely short, contributing to higher F2. This accounts for the wide separation of the two lowest dark horizontal bands in Figure 2.9, corresponding to F1 and F2, respectively. Rounding the lips has the effect of extending the front-of-tongue cavity, thus lowering F2. Typical values of F1 and F2 of American English vowels are listed in Table 2.5.

**Table 2.5** Phoneme labels and typical formant values for vowels of English.

| Vowel Labels | Mean F1 (Hz) | Mean F2 (Hz) |
|---|---|---|
| *iy (feel)* | 300 | 2300 |
| *ih (fill)* | 360 | 2100 |
| *ae (gas)* | 750 | 1750 |
| *aa (father)* | 680 | 1100 |
| *ah (cut)* | 720 | 1240 |
| *ao (dog)* | 600 | 900 |
| *ax (comply)* | 720 | 1240 |
| *eh (pet)* | 570 | 1970 |
| *er (turn)* | 580 | 1380 |
| *ow (tone)* | 600 | 900 |
| *uh (good)* | 380 | 950 |
| *uw (tool)* | 300 | 940 |

The characteristic F1 and F2 values for vowels are sometimes called formant targets, which are ideal locations for perception. Sometimes, due to fast speaking or other limitations on performance, the speaker cannot quite attain an ideal target before the articulators begin shifting to targets for the following phoneme, which is phonetic context dependent. Additionally, there is a special class of vowels that combine two distinct sets of F1/F2 targets.

These are called *diphthongs*. As the articulators move, the initial vowel targets glide smoothly to the final configuration. Since the articulators are working faster in production of a diphthong, sometimes the *ideal* formant target values of the component values are not fully attained. Typical diphthongs of American English are listed in Table 2.6.

**Table 2.6** The diphthongs of English.

| Diphthong Labels | Components |
|---|---|
| ay (t**ie)** | */aa/ ➔ /iy/* |
| ey (**a**te) | */eh/ ➔ /iy/* |
| oy (c**oi**n) | */ao/ ➔ /iy/* |
| aw (f**oul**) | */aa/ ➔ /uw/* |

Figure 2.17 shows the first two formants for a number of typical vowels.



**Figure 2.17** F1 and F2 values for articulations of some English vowels.

The major articulator for English vowels is the middle to rear portion of the tongue. The position of the tongue's surface is manipulated by large and powerful muscles in its root, which move it as a whole within the mouth. The linguistically important dimensions of movement are generally the ranges [front ⇔ back] and [high ⇔ low]. You can feel this movement easily. Say mentally, or whisper, the sound */iy/* (as in *see*) and then */aa/* (as in *father*). Do it repeatedly, and you will get a clear perception of the tongue movement from high to low. Now try */iy/* and then */uw/* (as in *blue*), repeating a few times. You will get a clear perception of place of articulation from front */iy/* to back */uw/*. Figure 2.18 shows a schematic characterization of English vowels in terms of relative tongue positions. There are two kinds of vowels: those in which tongue height is represented as a point and those in which it is represented as a vector.

Though the tongue hump is the major actor in vowel articulation, other articulators come into play as well. The most important secondary vowel mechanism for English and many other languages is lip rounding. Repeat the exercise above, moving from the */iy/* (*see*)

to the /uw/ (*blue*) position. Now rather than noticing the tongue movement, pay attention to your lip shape. When you say /iy/, your lips will be flat, slightly open, and somewhat spread. As you move to /uw/, they begin to *round out*, ending in a more puckered position. This lengthens the oral cavity during /uw/, and affects the spectrum in other ways.

**Figure 2.18** Relative tongue positions of English vowels [24].

Though there is always some controversy, linguistic study of phonetic abstractions, called *phonology*, has largely converged on the five binary features: +/- high, +/- low, +/- front, +/-back, and +/-round, plus the phonetically ambiguous but phonologically useful feature +/- tense, as adequate to uniquely characterize the major vowel distinctions of Standard English (and many other languages). Obviously, such a system is a little bit too free with logically contradictory specifications, such as [+high, +low], but these are excluded from real-world use. These features can be seen in Table 2.7.

**Table 2.7** Phonological (abstract) feature decomposition of basic English vowels.

| Vowel | high | low | front | back | round | tense |
|-------|------|-----|-------|------|-------|-------|
| iy | + | - | + | - | - | + |
| ih | + | - | + | - | - | - |
| ae | - | + | + | - | - | + |
| aa | - | + | - | - | - | + |
| ah | - | - | - | - | - | + |
| ao | - | + | - | + | + | + |
| ax | - | - | - | - | - | - |
| eh | - | - | + | - | - | - |
| ow | - | - | - | + | + | + |
| uh | + | - | - | + | - | - |
| uw | + | - | - | + | - | + |

This kind of abstract analysis allows researchers to make convenient statements about classes of vowels that behave similarly under certain conditions. For example, one may speak simply of the high vowels to indicate the set /*iy, ih, uh, uw*/.

## 2.2.1.2.    Consonants

Consonants, as opposed to vowels, are characterized by significant constriction or obstruction in the pharyngeal and/or oral cavities. Some consonants are voiced; others are not. Many consonants occur in pairs, that is, they share the same configuration of articulators, and one member of the pair additionally has voicing which the other lacks. One such pair is /*s, z*/, and the voicing property that distinguishes them shows up in the non-periodic noise of the initial segment /*s*/ in Figure 2.5 as opposed to the voiced consonant end-phone, /*z*/. Manner of articulation refers to the articulation mechanism of a consonant. The major distinctions in manner of articulation are listed in Table 2.8.

**Table 2.8** Consonant manner of articulation.

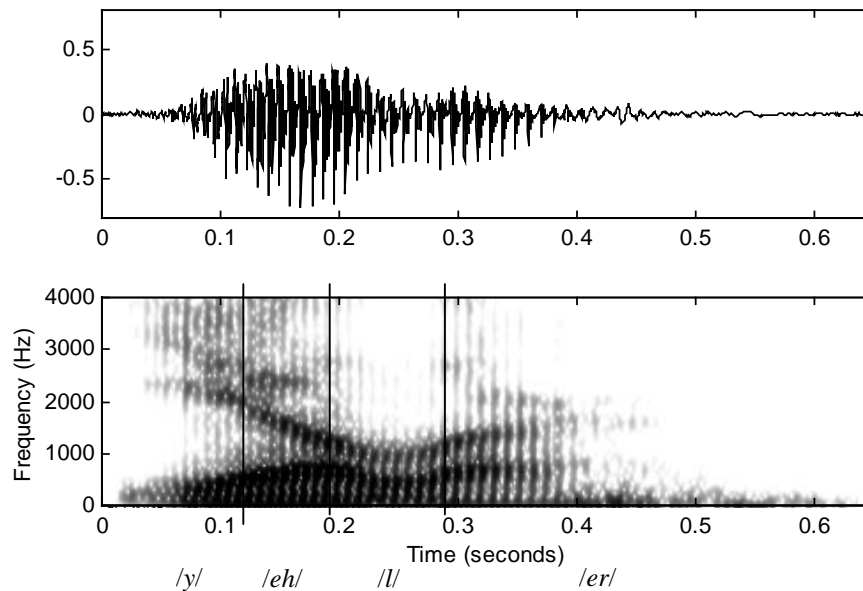| **Manner** | **Sample Phone** | **Example Words** | **Mechanism** |
|---|---|---|---|
| Plosive | /p/ | tat, tap | Closure in oral cavity |
| Nasal | /m/ | team, meet | Closure of nasal cavity |
| Fricative | /s/ | sick, kiss | Turbulent airstream noise |
| Retroflex liquid | /r/ | rat, tar | Vowel-like, tongue high and curled back |
| Lateral liquid | /l/ | lean, kneel | Vowel-like, tongue central, side airstream |
| Glide | /y/,/w/ | yes, well | Vowel-like |

The English phones that typically have voicing without complete obstruction or narrowing of the vocal tract are called *semivowels* and include /*l, r*/, the *liquid* group, and /*y, w*/, the *glide* group. Liquids, glides, and vowels are all *sonorant*, meaning they have continuous voicing. Liquids /*l*/ and /*r*/ are quite vowel-like and in fact may become *syllabic* or act entirely as vowels in certain positions, such as the *l* at the end of *edible*. In /*l*/, the airstream flows around the sides of the tongue, leading to the descriptive term *lateral*. In /*r*/, the tip of the tongue is curled back slightly, leading to the descriptive term *retroflex*. Figure 2.19 shows some semivowels.

Glides /*y, w*/ are basically vowels /*iy, uw*/ whose initial position within the syllable require them to be a little shorter and to lack the ability to be stressed, rendering them just different enough from true vowels that they are classed as a special category of consonant. Pre-vocalic glides that share the syllable-initial position with another consonant, such as the /*y*/ in the second syllable of *computer /k uh m . p **y** uw . t er/*, or the /*w*/ in *quick /k **w** ih k/*, are sometimes called *on-glides*. The semivowels, as a class, are sometimes called *approximants,* meaning that the tongue approaches the top of the oral cavity, but does not completely contact so as to obstruct the air flow.

Even the non-sonorant consonants that require complete or close-to-complete obstruction may still maintain some voicing before or during the obstruction, until the pressure differential across the glottis starts to disappear, due to the closure. Such voiced consonants

include /b,d,g, z, zh, v/. They have a set of counterparts that differ only in their characteristic lack of voicing: /p,t,k, s, sh, f/.



**Figure 2.19** Spectrogram for the word *yeller*, showing semivowels /y/, /l/, /er/ (approximate phone boundaries shown with vertical lines).

Nasal consonants /m,n/ are a mixed bag: the oral cavity has significant constriction (by the tongue or lips), yet the voicing is continuous, like that of the sonorants, because, with the velar flap open, air passes freely through the nasal cavity, maintaining a pressure differential across the glottis.

A consonant that involves complete blockage of the oral cavity is called an obstruent stop, or plosive consonant. These may be voiced throughout if the trans-glottal pressure drop can be maintained long enough, perhaps through expansion of the wall of the oral cavity. In any case, there can be voicing for the early sections of stops. Voiced, unvoiced pairs of stops include: /b,p/, /d,t/, and /g,k/. In viewing the waveform of a stop, a period of silence corresponding to the oral closure can generally be observed. When the closure is removed (by opening the constrictor, which may be lips or tongue), the trapped air rushes out in a more or less sudden manner. When the upper oral cavity is unimpeded, the closure of the vocal folds themselves can act as the initial blocking mechanism for a type of stop heard at the very beginning of vowel articulation in vowel-initial words like *atrophy*. This is called a *glottal stop*. Voiceless plosive consonants in particular exhibit a characteristic aperiodic *burst* of energy at the (articulatory) point of closure as shown in Figure 2.20 just prior to /i/. By comparison, the voicing of voiced plosive consonants may not always be obvious in a spectrogram.
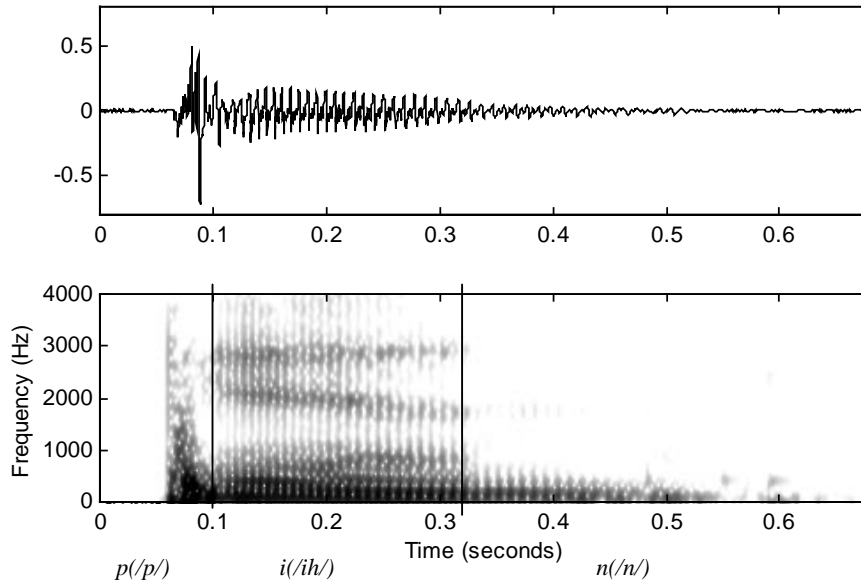
**Figure 2.20** Spectrogram: stop release *burst* of /*p*/ in the word *pin*.

A consonant that involves nearly complete blockage of some position in the oral cavity creates a narrow stream of turbulent air. The friction of this air stream creates a non-periodic hiss-like effect. Sounds with this property are called fricatives and include /s, z/. There is no voicing during the production of *s*, while there can be voicing (in addition to the frication noise), during the production of *z*, as discussed above. /*s, z*/ have a common place of articulation, as explained below, and thus form a natural similarity class. Though controversial, /*h*/ can also be thought of as a (glottal) fricative. /*s*/ in word-initial position and /*z*/ in word-final position are exemplified in Figure 2.5.

Some sounds are complex combinations of manners of articulation. For example, the *affricates* consist of a stop (e.g., /*t*/), followed by a fricative [e.g., /*sh*/) combining to make a unified sound with rapid phases of closure and continuancy (e.g., {*t* + *sh*) = *ch* as in *church*). The affricates in English are the voiced/unvoiced pairs: /*j*/ (*d* + *zh*) and /*ch*/ (*t* + *sh*). The complete consonant inventory of English is shown in Table 2.9.

Consider the set /*m*/, /*n*/, /*ng*/ from Table 2.9. They are all voiced nasal consonants, yet they sound distinct to us. The difference lies in the location of the major constriction along the top of the oral cavity (from lips to velar area) that gives each consonant its unique quality. The articulator used to touch or approximate the given location is usually some spot along the length of the tongue. As shown in Figure 2.21, the combination of articulator and place of articulation gives each consonant its characteristic sound:
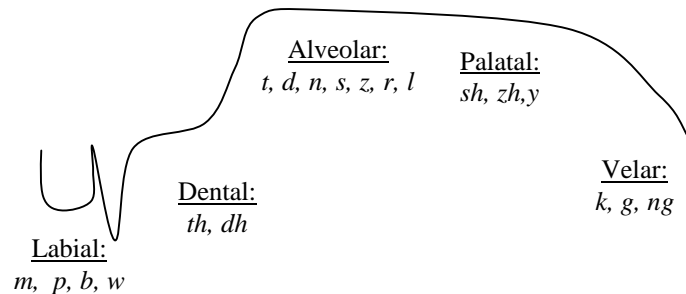
The *labial* consonants have their major constriction at the lips. This includes /*p*/, /*b*/ (these two differ only by manner of articulation) and /*m*/ and /*w*/.

The class of *dental or labio-dental* consonants includes /f, v/ and /th, dh/ (the members of these groups differ in manner, not place).

*Alveolar* consonants bring the front part of the tongue, called the tip or the part behind the tip called the blade, into contact or approximation to the alveolar ridge, rising semi-vertically above and behind the teeth. These include /t, d, n, s, z, r, l/. The members of this set again differ in manner of articulation (voicing, continuity, nasality), rather than place.

*Palatal* consonants have approximation or constriction on or near the roof of the mouth, called the palate. The members include /sh, zh, y/.

*Velar* consonants bring the articulator (generally the back of the tongue), up to the rearmost top area of the oral cavity, near the velar flap. Velar consonants in English include /k, g/ (differing by voicing) and the nasal continuant /ng/.

**Table 2.9** Manner of articulation of English consonants.

| Consonant Labels | Consonant Examples | Voiced? | Manner |
|:---:|:---:|:---:|:---:|
| *b* | **b**ig, a**b**le, ta**b** | + | plosive |
| *p* | **p**ut, o**p**en, ta**p** | - | plosive |
| *d* | **d**ig, i**d**ea, wa**d** | + | plosive |
| *t* | **t**alk, sa**t** | - | plosive |
| *g* | **g**ut, an**g**le, ta**g** | + | plosive |
| *k* | **c**ut, oa**k**en, ta**k**e | - | plosive |
| *v* | **v**at, o**v**er, ha**v**e | + | fricative |
| *f* | **f**ork, a**f**ter, i**f** | - | fricative |
| *z* | **z**ap, la**z**y, ha**z**e | + | fricative |
| *s* | **s**it, ca**s**t, to**ss** | - | fricative |
| *dh* | **th**en, fa**th**er, scy**th**e | + | fricative |
| *th* | **th**in, no**th**ing, tru**th** | - | fricative |
| *zh* | **g**enre, a**z**ure, bei**g**e | + | fricative |
| *sh* | **sh**e, cu**sh**ion, wa**sh** | - | fricative |
| *jh* | **j**oy, a**g**ile, e**dge** | + | affricate |
| *ch* | **ch**in, ar**ch**er, mar**ch** | - | affricate |
| *l* | **l**id, e**l**bow, sai**l** | + | lateral |
| *r* | **r**ed, pa**r**t, fa**r** | + | retroflex |
| *y* | **y**acht, on**i**on, **y**ard | + | glide |
| *w* | **w**ith, a**w**ay | + | glide |
| *hh* | **h**elp, a**h**ead, **h**otel | + | fricative |
| *m* | **m**at, a**m**id, ai**m** | + | nasal |
| *n* | **n**o, e**n**d, pa**n** | + | nasal |
| *ng* | si**ng**, a**ng**er, dri**n**k | + | nasal |

With the place terminology, we can complete the descriptive inventory of English consonants, arranged by manner (rows), place (columns) and voiceless/voiced (pairs in cells) as illustrated in Table 2.10.

**Figure 2.21** The major places of consonant articulation with respect to human mouth.

**Table 2.10** The consonants of English arranged by place (columns) and manner (rows).

|  | **Labial** | **Labio-dental** | **Dental** | **Alveolar** | **Palatal** | **Velar** | **Glottal** |
|---|---|---|---|---|---|---|---|
| Plosive | *p b* |  |  | *t d* |  | *k g* | *?* |
| Nasal | *m* |  |  | *n* |  | *ng* |  |
| Fricative |  | *f v* | *th dh* | *s z* | *sh zh* |  | *h* |
| Retroflex Sonorant |  |  |  | *r* |  |  |  |
| Lateral sonorant |  |  |  | *l* |  |  |  |
| Glide | *w* |  |  |  | *y* |  |  |

## 2.2.1.3.    Phonetic Typology

The oral, nasal, pharyngeal, and glottal mechanisms actually make available a much wider range of effects than English happens to use. So, it is expected that other languages would utilize other vocal mechanisms, in an internally consistent but essentially arbitrary fashion, to represent their lexicons. In addition, often a vocal effect that is part of the systematic linguistic phonetics of one language is present in others in a less codified, but still perceptible, form. For example, Japanese vowels have a characteristic distinction of length that can be hard for non-natives to perceive and to use when learning the language. The words *kado* (*corner*) and *kaado* (*card*) are spectrally identical, differing only in that *kado* is much shorter in all contexts. The existence of such minimally-contrasting pairs is taken as conclusive evidence that *length* is phonemically distinctive for Japanese. As noted above, what is linguistically distinctive in any one language is generally present as a less *meaningful* signaling dimension in other languages. Thus, vowel length can be manipulated in any English word as well, but this occurs either consciously for emphasis or humorous effect, or unconsciously and very predictably at clause and sentence end positions, rather than to signal lexical identity in all contexts, as in Japanese.

Other interesting sounds that the English language makes no linguistic use of include the trilled *r* sound and the implosive. The trilled *r* sound is found in Spanish, distinguishing

(for example) the words *pero* (*but*) and *perro* (*dog*). This trill could be found in times past as a non-lexical sound used for emphasis and interest by American circus ringmasters and other showpersons.

While the world's languages have all the variety of manner of articulation exemplified above and a great deal more, the primary dimension lacking in English that is exploited by a large subset of the world's languages is pitch variation. Many of the huge language families of Asia and Africa are tonal, including all varieties of Chinese. A large number of other languages are not considered strictly tonal by linguistics, yet they make systematic use of pitch contrasts. These include Japanese and Swedish. To be considered tonal, a language should have lexical meaning contrasts cued by pitch, just as the lexical meaning contrast between *pig* and *big* is cued by a voicing distinction in English. For example, Mandarin Chinese has four primary tones (tones can have minor context-dependent variants just like ordinary phones, as well):

**Table 2.11** The contrastive tones of Mandarin Chinese.

| Tone | Shape | Example | Chinese | Meaning |
|---|---|---|---|---|
| 1 | High level | *ma* | 妈 | mother |
| 2 | High rising | *ma* | 麻 | numb |
| 3 | Low rising | *ma* | 马 | horse |
| 4 | High falling | *ma* | 骂 | to scold |

Though English does not make systematic use of pitch in its inventory of word contrasts, nevertheless, as we always see with any possible phonetic effect, pitch is systematically varied in English to signal a speaker's emotions, intentions, and attitudes, and it has some linguistic function in signaling grammatical structure as well. Pitch variation in English will be considered in more detail in Chapter 15.
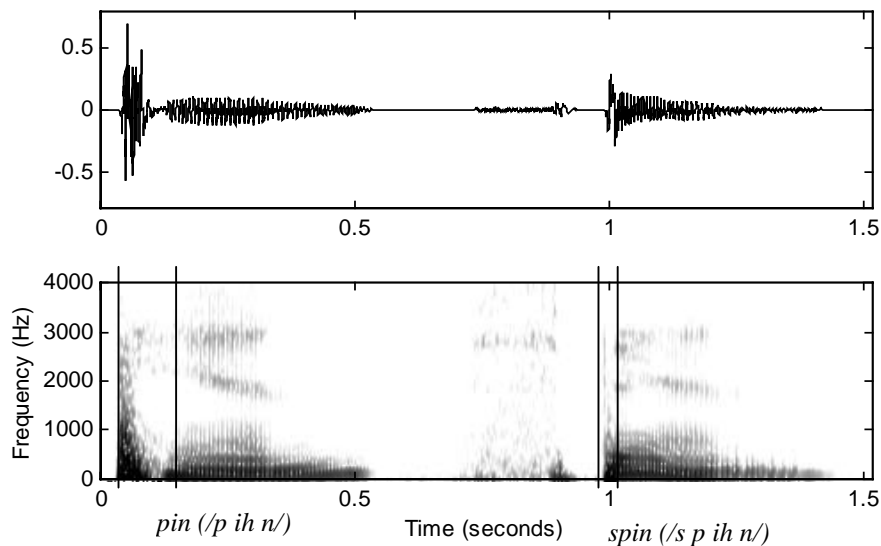
## 2.2.2.    The Allophone: Sound and Context

The vowel and consonant charts provide abstract symbols for the phonemes - major sound distinctions. Phonemic units should be correlated with potential meaning distinctions. For example, the change created by holding the tongue high and front (/*iy*/) vs. directly down from the (frontal) position for /*eh*/, in the consonant context /*m _ n*/, corresponds to an important meaning distinction in the lexicon of English: *mean* /*m iy n*/ vs. *men* /*m eh n*/. This meaning contrast, conditioned by a pair of rather similar sounds, in an identical context, justifies the inclusion of /*iy*/ and /*eh*/ as logically separate distinctions.

However, one of the fundamental, meaning-distinguishing sounds is often modified in some systematic way by its phonetic neighbors. The process by which neighboring sounds influence one another is called *coarticulation*. Sometimes, when the variations resulting from coarticulatory processes can be consciously perceived, the modified phonemes are called *allophones*. Allophonic differences are always *categorical*, that is, they can be understood and denoted by means of a small, bounded number of symbols or diacritics on the basic phoneme symbols.

As an experiment, say the word *like* to yourself. Feel the front of the tongue touching the alveolar ridge (cf. Figure 2.21) when realizing the initial phoneme /l/. This is one allophone of /l/, the so-called *light* or *clear* /l/. Now say *kill*. In this word, most English speakers will no longer feel the front part of the tongue touch the alveolar ridge. Rather, the /l/ is realized by stiffening the broad mid section of the tongue in the rear part of the mouth while the continuant airstream escapes laterally. This is another allophone of /l/, conditioned by its syllable-final position, called the *dark* /l/. Predictable contextual effects on the realization of phones can be viewed as a nuisance for speech recognition, as will be discussed in Chapter 9. On the other hand, such variation, because it is systematic, could also serve as a cue to the syllable, word, and prosodic structure of speech.

Now experiment with the sound /p/ by holding a piece of tissue in front of your mouth while saying the word *pin* in a normal voice. Now repeat this experiment with *spin*. For most English speakers, the word *pin* produces a noticeable puff of air, called aspiration. But the same phoneme, /p/, embedded in the consonant cluster /sp/ loses its aspiration (burst, see the lines bracketing the /p/ release in *pin* and *spin* in Figure 2.22), and because these two types of /p/ are in complementary distribution (completely determined by phonetic and syllabic context), the difference is considered allophonic.



**Figure 2.22** Spectrogram: bursts of *pin* and *spin*. The relative duration of a p-burst in different phonetic contexts is shown by the differing width of the area between the vertical lines.

Try to speak the word *bat* in a framing phrase *say bat again*. Now speak *say bad again*. Can you feel the length difference in the vowel /ae/? A vowel before a voiced consonant e.g., /d/, seems typically longer than the same vowel before the unvoiced counterpart, in this case /t/.

A sound phonemicized as /t/ or /d/, that is, a stop made with the front part of the tongue, may be reduced to a quick tongue tap that has a different sound than either /t/ or /d/ in fuller contexts. This process is called flapping. It occurs when /t/ or /d/ closes a stressed vowel (coda position) followed by an unstressed vowel, as in: *bitter, batter, murder, quarter, humidity,* and can even occur across words as long as the preconditions are met, as in *you can say that again.* Sometimes the velar flap opens too soon (anticipation), giving a characteristically nasal quality to some pre-nasal vowels such as /ae/ in *ham* vs. *had*. We have a more detailed discussion on allophones in Chapter 9.

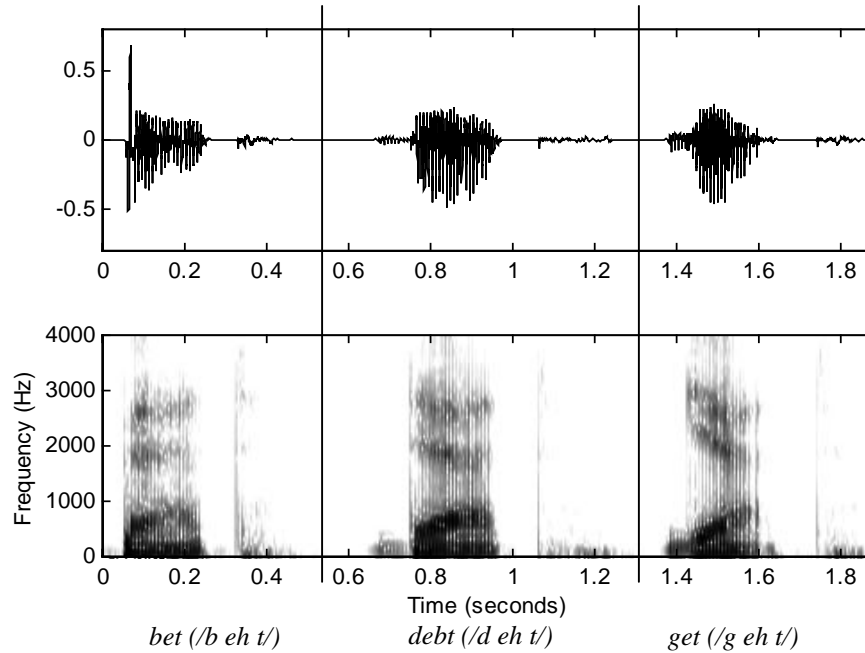## 2.2.3. Speech Rate and Coarticulation

In addition to allophones, there are other variations in speech for which no small set of established categories of variation can be established. These are *gradient*, existing along a scale for each relevant dimension, with speakers scattered widely. In general, it is harder to become consciously aware of coarticulation effects than of allophonic alternatives.

Individual speakers may vary their rates according to the content and setting of their speech, and there may be great inter-speaker differences as well. Some speakers may pause between every word, while others may speak hundreds of words per minute with barely a pause between sentences. At the faster rates, formant targets are less likely to be fully achieved. In addition, individual allophones may merge.

For example [20], consider the utterance *Did you hit it to Tom*? The pronunciation of this utterance is */d ih d y uw h ih t t uw t aa m/*. However, a realistic, casual rendition of this sentence would appear as */d ih jh ax hh ih dx ih t ix t aa m/,* where /ix/ is a reduced schwa /ax/ that is short and often unvoiced, and /dx/ is a kind of shortened, indistinct stop, intermediate between /d/ and /t/. The following five phonologic rules have operated on altering the pronunciation in the example:

> Palatalization of /d/ before /y/ in *di_d_ you*
>
> Reduction of unstressed /u/ to schwa in *y_ou_*
>
> Flapping of intervocalic /t/ in *hi_t_ it*
>
> Reduction of schwa and devoicing of /u/ in *t_o_*
>
> Reduction of geminate (double consonant) /t/ in *i_t_ _to_*

There are also coarticulatory influences in the spectral appearance of speech sounds, which can only be understood at the level of spectral analysis. For example, in vowels, consonant neighbors can have a big effect on formant trajectories near the boundary. Consider the differences in *F1* and *F2* in the vowel /eh/ as realized in words with different initial consonants *bet*, *debt*, and *get*, corresponding to the three major places of articulation (labial, alveolar, and velar), illustrated in Figure 2.23. You can see the different relative spreads of *F1* and *F2* following the initial stop consonants.

**Figure 2.23** Spectrogram*: bet, debt*, and *get* (separated by vertical lines). Note different relative spreads of *F1* and *F2* following the initial stop consonants in each word.

Now let's see different consonants following the same vowel, *ebb*, *head*, and *egg*. In Figure 2.23, the coarticulatory effect is *perseverance*; i.e., in the early part of the vowel the articulators are still somewhat set from realization of the initial consonant. In the *ebb*, *head*, and *egg* examples shown in Figure 2.24, the coarticulatory effect is *anticipation*; i.e., in the latter part of the vowel the articulators are moving to prepare for the upcoming consonant articulation. You can see the increasing relative spread of *F1* and *F2* at the final vowel-consonant transition in each word.

## 2.3.    SYLLABLES AND WORDS

Phonemes are small building blocks. To contribute to language meaning, they must be organized into longer cohesive spans, and the units so formed must be combined in characteristic patterns to be meaningful, such as syllables and words in the English language.
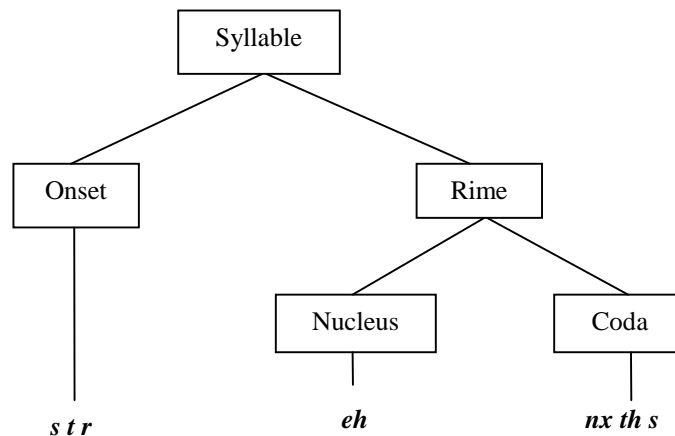
**Figure 2.24** Spectrogram: *ebb*, *head*, and *egg*. Note the increasing relative spread of *F1* and *F2* at the final vowel-consonant transition in each word.

## 2.3.1.    Syllables

An intermediate unit, the *syllable*, is sometimes thought to interpose between the phones and the word level. The syllable is a slippery concept, with implications for both production and perception. Here we will treat it as a perceptual unit. Syllables are generally centered around vowels in English, giving two perceived syllables in a word like *tomcat*: /tOm-cAt/. To completely parse a word into syllables requires making judgments of consonant affiliation (with the syllable peak vowels). The question of whether such judgments should be based on articulatory or perceptual criteria, and how they can be rigorously applied, remains unresolved.

Syllable centers can be thought of as *peaks* in sonority (high-amplitude, periodic sections of the speech waveform). These sonority peaks have affiliated *shoulders* of strictly non-increasing sonority. A scale of sonority can be used, ranking consonants along a continuum of stops, affricates, fricatives, and approximants. So, in a word like *verbal*, the syllabification would be *ver-bal*, or *verb-al*, but not *ve-rbal*, because putting the approximant /r/ before the stop /b/ in the second syllable would violate the non-decreasing sonority requirement heading into the syllable.

As long as the sonority conditions are met, the exact affiliation of a given consonant that could theoretically affiliate on either side can be ambiguous, unless determined by higher-order considerations of word structure, which may block affiliation. For example, in a word like *beekeeper*, an abstract boundary in the compound between the component words *bee* and *keeper* keeps us from accepting the syllable parse: *beek-eeper*, based on lexical interpretation. However, the same phonetic sequence in *beaker* could, depending on one's theory of syllabicity, permit affiliation of the *k: beak-er*. In general, the syllable is a unit that has intuitive plausibility but remains difficult to pin down precisely.

```
                        ┌──────────┐
                        │ Syllable │
                        └──────────┘
                       ╱            ╲
              ┌────────┐            ┌──────┐
              │ Onset  │            │ Rime │
              └────────┘            └──────┘
                  │                 ╱        ╲
                  │          ┌─────────┐   ┌──────┐
                  │          │ Nucleus │   │ Coda │
                  │          └─────────┘   └──────┘
                  │               │            │
                 s t r           eh         nx th s
```

**Figure 2.25** The word/syllable *strengths* (*/s t r eh nx th s/*) is a longest syllable of English.

Syllables are thought (by linguistic theorists) to have internal structure, and the terms used are worth knowing. Consider a big syllable such as *strengths /s t r eh nx th s/*. This consists of a vowel peak, called the *nucleus*, surrounded by the other sounds in characteristic positions. The *onset* consists of initial consonants if any, and the rime is the nucleus with trailing consonants (the part of the syllable that matters in determining poetic rhyme). The coda consists of consonants in the rime following the nucleus (in some treatments, the last consonant in a final cluster would belong to an *appendix*). This can be diagrammed as a syllable parse tree as shown in Figure 2.25. The syllable is sometimes thought to be the primary domain of coarticulation, that is, sounds within a syllable influence one another's realization more than the same sounds separated by a syllable boundary.

## 2.3.2.    Words

The concept of words seems intuitively obvious to most speakers of Indo-European languages. It can be loosely defined as a lexical item, with an agreed-upon meaning in a given speech community, that has the freedom of syntactic combination allowed by its type (noun, verb, etc.).

In spoken language, there is a segmentation problem: words *run together* unless affected by a disfluency (unintended speech production problem) or by the deliberate placement of a pause (silence) for some structural or communicative reason. This is surprising to many people, because literacy has conditioned speakers/readers of Indo-European languages to expect a *blank space* between words on the printed page. But in speech, only a few true pauses (the aural equivalent of a blank space) may be present. So, what appears to the reading eye as *never give all the heart, for love* would appear to the ear, if we simply use letters to stand for their corresponding English speech sounds, as *nevergivealltheheart   forlove* or, in phonemes, as *n eh v er g ih v ah l dh ax h aa r t \\ f ao r l ah v.* The \\ symbol marks a linguistically motivated pause, and the units so formed are sometimes called *intonation phrases*, as explained in Chapter 15.

Certain facts about word structure and combinatorial possibilities are evident to most native speakers and have been confirmed by decades of linguistic research. Some of these facts describe relations among words when considered in isolation, or concern groups of related words that seem intuitively similar along some dimension of form or meaning - these properties are *paradigmatic*. Paradigmatic properties of words include part-of-speech, inflectional and derivational morphology, and compound structure. Other properties of words concern their behavior and distribution when combined for communicative purposes in fully functioning utterances – these properties are *syntagmatic*.

### 2.3.2.1.    Lexical Part-of-Speech

Lexical part-of-speech (POS) is a primitive form of linguistic theory that posits a restricted inventory of word-type categories, which capture generalizations of word forms and distributions. Assignment of a given POS specification to a word is a way of summarizing certain facts about its potential for syntagmatic combination. Additionally, paradigms of word formation processes are often similar within POS types and subtypes as well. The word properties upon which POS category assignments are based may include affixation behavior, very abstract semantic typologies, distributional patterns, compounding behavior, historical development, productivity and generalizabilty, and others.

A typical set of POS categories would include *noun*, *verb*, *adjective*, *adverb, interjection*, *conjunction*, *determiner*, *preposition*, and *pronoun*. Of these, we can observe that certain classes of words consist of infinitely large membership. This means new members can be added at any time. For example, the category of noun is constantly expanded to accommodate new inventions, such as Velcro or Spandex. New individuals are constantly being born, and their names are a type of noun called *proper noun*. The proliferation of words using the descriptive prefix *cyber* is another recent set of examples: *cyberscofflaw, cybersex*, and even *cyberia* illustrate the infinite creativity of humans in manipulating word structure to express new shades of meaning, frequently by analogy with, and using fragments of, existing vocabulary. Another example is the neologism *sheeple*, a noun combining the forms and meanings of *sheep* and *people* to refer to large masses of people who lack the capacity or willingness to take independent action. We can create new words whenever we like, but they had best fall within the predictable paradigmatic and syntagmatic patterns of use summarized by the existing POS generalizations, or there will be little hope of their adoption by

any other speaker. These open POS categories are listed in Table 2.12. Nouns are inherently referential. They refer to persons, places, and things. Verbs are predicative; they indicate relations between entities and properties of entities, including participation in events. Adjectives typically describe and more completely specify noun reference, while adverbs describe, intensify and more completely specify verbal relations. Open-class words are sometimes called *content* words, for their referential properties.

In contrast to the open-class categories, certain other categories of words only rarely and very slowly admit new members over the history of English development. These closed POS categories are shown in Table 2.13. The closed-category words are fairly stable over time. Conjunctions are used to join larger syntactically complete phrases. Determiners help to narrow noun reference possibilities. Prepositions denote common spatial and temporal relations of objects and actions to one another. Pronouns provide a convenient substitute for noun phrases that are fully understood from context. These words denote grammatical relations of other words to one another and fundamental properties of the world and how humans understand it. They can, of course, change slowly; for example, the Middle English pronoun *thee* is no longer in common use. The closed-class words are sometimes called *function* words.

**Table 2.12** Open POS categories.

| Tag | Description | Function | Example |
|---|---|---|---|
| N | Noun | Names entity | *cat* |
| V | Verb | Names event or condition | *forget* |
| Adj | Adjective | Descriptive | *yellow* |
| Adv | Adverb | Manner of action | *quickly* |
| Interj | Interjection | Reaction | *oh!* |

**Table 2.13** Closed POS categories.

| Tag | Description | Function | Example |
|---|---|---|---|
| Conj | Conjunction | Coordinates phrases | *and* |
| Det | Determiner | Indicates definiteness | *the* |
| Prep | Preposition | Relations of time, space, direction | *from* |
| Pron | Pronoun | Simplified reference | *she* |

The set of POS categories can be extended indefinitely. Examples can be drawn from the Penn Treebank project (http://www.cis.upenn.edu/ldc) as shown in Table 2.14, where you can find the proliferation of sub-categories, such as *Verb, base form* and *Verb, past tense*. These categories incorporate *morphological* attributes of words into the POS label system discussed in Section 2.3.2.2.

POS tagging is the process of assigning a part-of-speech or other lexical class marker to each word in a corpus. There are many algorithms exist to automatically tag input sentences into a set of tags. Rule-based methods [45], hidden Markov models (see Chapter 8) [23, 29, 46], and machine-learning methods [6] are used for this purpose.

### 2.3.2.2.    Morphology

Morphology is about the subparts of words, i.e., the patterns of word formation including inflection, derivation, and the formation of compounds. English mainly uses prefixes and suffixes to express *inflection* and *derivational* morphology.

**Table 2.14** Treebank POS categories – an expanded inventory.

| String | Description | Example |
|--------|-------------|---------|
| CC | Coordinating conjunction | and |
| CD | Cardinal number | two |
| DT | Determiner | the |
| EX | Existential *there* | there (*There was an old lady*) |
| FW | Foreign word | *omerta* |
| IN | Preposition, subord. conjunction | over, but |
| JJ | Adjective | yellow |
| JJR | Adjective, comparative | better |
| JJS | Adjective, superlative | best |
| LS | List item marker | |
| MD | Modal | might |
| NN | Noun, singular or mass | rock, water |
| NNS | Noun, plural | rocks |
| NNP | Proper noun, singular | Joe |
| NNPS | Proper noun, plural | Red Guards |
| PDT | Predeterminer | all (*all the girls*) |
| POS | Possessive ending | 's |
| PRP | Personal pronoun | I |
| PRP$ | Possessive pronoun | mine |
| RB | Adverb | quickly |
| RBR | Adverb, comparative | higher (*shares closed higher.*) |
| RBS | Adverb, superlative | highest (*he jumped highest of all.*) |
| RP | Particle | up ( *take up the cause*) |
| TO | *to* | *to* |
| UH | Interjection | hey! |
| VB | Verb, base form | choose |
| VBD | Verb, past tense | chose |
| VBG | Verb, gerund or present participle | choosing |
| VBN | Verb, past participle | chosen |
| VBP | Verb, non-third person sing. present | jump |
| VBZ | Verb, third person singular present | jumps |
| WDT | Wh-determiner | which |
| WP | Wh-pronoun | who |
| WP$ | Possessive wh-pronoun | whose |
| WRB | Wh-adverb | when (*When he came, it was late.*) |

*Inflectional morphology* deals with variations in word form that reflect the contextual situation of a word in phrase or sentence syntax, and that rarely have direct effect on interpretation of the fundamental meaning expressed by the word. English inflectional morphology is relatively simple and includes person and number agreement and tense markings only. The variation in *cats* (vs. *cat*) is an example. The plural form is used to refer to an indefinite number of cats greater than one, depending on a particular situation. But the basic POS category (*noun*) and the basic meaning (*felis domesticus*) are not substantially affected. Words related to a common lemma via inflectional morphology are said to belong to a common paradigm, with a single POS category assignment. In English, common paradigm types include the verbal set of affixes (pieces of words): *-s, -ed, -ing*, the noun set: *-s*, and the adjectival *-er, -est*. Note that sometimes the base form may change spelling under affixation, complicating the job of automatic textual analysis methods. For historical reasons, certain paradigms may consist of highly idiosyncratic irregular variation as well, e.g., *go, going, went, gone* or *child, children*. Furthermore, some words may belong to defective paradigms, where only the singular (noun: *equipment*) or the plural (noun: *scissors*) is provided for.

In *derivational morphology*, a given root word may serve as the source for wholly new words, often with POS changes as illustrated in Table 2.15. For example, the terms *racial* and *racist*, though presumably based on a single root word *race,* have different POS possibilities (*adjective* vs. *noun-adjective*) and meanings. Derivational processes may induce pronunciation change or stress shift (e.g., *el**e**ctric* vs. *electr**i**city*). In English, typical derivational affixes (pieces of words) that are highly productive include prefixes and suffixes: *re-, pre-, -ial, -ism, -ish, -ity, -tion, -ness, -ment, -ious, -ify, -ize*, and others. In many cases, these can be added successively to create a complex layered form.

**Table 2.15** Examples of stems and their related forms across POS categories.

| Noun | Verb | Adjective | Adverb |
|------|------|-----------|--------|
| *criticism* | *criticize* | *critical* | *critically* |
| *fool* | *fool* | *foolish* | *foolishly* |
| *industry, industrialization* | *industrialize* | *industrial,industrious* | *industriously* |
| *employ, employee, employer* | *employ* | *employable* | *employably* |
| *certification* | *certify* | *certifiable* | *certifiably* |

Generally, word formation operates in layers, according to a kind of word syntax: *(deriv-prefix)\* **root** (root)\* (deriv-suffix)\* (infl-suffix).* This means that one or more *roots* can be compounded in the inner layer, with one or more optional *derivational prefixes*, followed by any number of optional *derivational suffixes*, capped off with no more than one *inflectional suffix*. There are, of course, limits on word formation, deriving both from semantics of the component words and simple lack of imagination. An example of a nearly maximal word in English might be *autocyberconceptualizations*, meaning (perhaps!) multiple instances of automatically creating computer-related concepts. This word lacks only compounding to be truly maximal. This word has a derivational prefix *auto-*, two root forms compounded (*cyber*

and *concept*, though some may prefer to analyze *cyber-* as a prefix), three derivational suffixes (*-ual, ize, -ation*), and is capped off with the plural inflectional suffix for nouns, *-s.*

### 2.3.2.3. Word Classes

POS classes are based on traditional grammatical and lexical analysis. With improved computational resources, it has become possible to examine words in context and assign words to groups according to their actual behavior in real text and speech from a statistical point of view. These kinds of classifications can be used in language modeling experiments for speech recognition, text analysis for text-to-speech synthesis, and other purposes.

One of the main advantages of word classification is its potential to derive more refined classes than traditional POS, while only rarely actually crossing traditional POS group boundaries. Such a system may group words automatically according to the similarity of usage with respect to their word neighbors. Consider classes automatically found by the classification algorithms of Brown *et al*. [7]:

{Friday Monday Thursday Wednesday Tuesday Saturday Sunday weekends}
{great big vast sudden mere sheer gigantic lifelong scant colossal}
{down backwards ashore sideways southward northward overboard aloft adrift}
{mother wife father son husband brother daughter sister boss uncle}
{John George James Bob Robert Paul William Jim David Mike}
{feet miles pounds degrees inches barrels tons acres meters bytes}

You can see that words are grouped together based on the semantic meaning, which is different from word classes created purely from syntactic point of view. Other types of classification are also possible, some of which can identify semantic relatedness across traditional POS categories. Some of the groups derived from this approach may include follows:

{problems problem solution solve analyzed solved solving}
{write writes writing written wrote pen}
{question questions asking answer answers answering}
{published publication author publish writer titled}

## 2.4. SYNTAX AND SEMANTICS

Syntax is the study of the patterns of formation of sentences and phrases from words and the rules for the formation of grammatical sentences. Semantics is another branch of linguistics dealing with the study of meaning, including the ways meaning is structured in language and changes in meaning and form over time.

## 2.4.1.     Syntactic Constituents

Constituents represent the way a sentence can be divided into its grammatical subparts as constrained by common grammatical patterns (which implicitly incorporate normative judgments on acceptability). Syntactic constituents at least respect, and at best explain, the linear order of words in utterances and text. In this discussion, we will not strictly follow any of the many theories of syntax but will instead bring out a few basic ideas common to many approaches. We will not attempt anything like a complete presentation of the grammar of English but instead focus on a few simple phenomena.

Most work in syntactic theory has adopted machinery from traditional grammatical work on written language. Rather than analyze toy sentences, let's consider what kinds of superficial syntactic patterns are lurking in a random chunk of serious English text, excerpted from David Thoreau's essay *Civil Disobedience* [43]:

*The authority of government, even such as I am willing to submit to - for I will cheerfully obey those who know and can do better than I, and in many things even those who neither know nor can do so well - is still an impure one: to be strictly just, it must have the sanction and consent of the governed. It can have no pure right over my person and property but what I concede to it. The progress from an absolute to a limited monarchy, from a limited monarchy to a democracy, is a progress toward a true respect for the individual.*

### 2.4.1.1.     Phrase Schemata

Words may be combined to form phrases that have internal structure and unity. We use generalized schemata to describe the phrase structure. The goal is to create a simple, uniform template that is independent of POS category.

Let's first consider nouns, a fundamental category referring to persons, places, and things in the world. The noun and its immediate modifiers form a constituent called the noun phrase (NP). To generalize this, we consider a word of arbitrary category, say category X (which could be a noun N or a verb V.). The generalized rule for a phrase XP is *XP* $\Rightarrow$ *(modifiers) X-head (post-modifiers),* where *X* is the head, since it dominates the configuration and names the phrase. Elements preceding the head in its phrase are *premodifiers* and elements following the head are *postmodifiers*. *XP*, the culminating phrase node, is called a *maximal projection* of category *X*. We call the whole structure an *x-template*. Maximal projections, *XP*, are the primary currency of basic syntactic processes. The post-modifiers are usually maximal projections (another head, with its own post-modifiers forming an *XP* on its own) and are sometimes termed *complements*, because they are often required by the lexical properties of the head for a complete meaning to be expressed (e.g. when *X* is a preposition or verb). Complements are typically noun phrases (*NP*), prepositional phrases (*PP*), verb phrases (*VP*), or sentence/clause (*S*), which make an essential contribution to the head's reference or meaning, and which the head requires for semantic completeness. Premodifiers are likely to be adverbs, adjectives, quantifiers, and determiners, i.e., words that help to specify the meaning of the head but may not be essential for completing the meaning. With minor variations, the *XP* template serves for most phrasal types, based on the POS of the head (*N, V, ADJ*, etc.).

For *NP*, we thus have *NP* $\Rightarrow$ *(det) (modifier)* **head-noun** *(post-modifier).* This rule describes an *NP* (noun phrase - left side of arrow) in terms of its optional and required internal contents (right side of the arrow). *Det* is a word like *the* or *a* that helps to resolve the reference to a specific or an unknown instance of the noun. The *modifier* gives further information about the noun. The *head* of the phrase, and the only mandatory element, is the noun itself. *Post-modifiers* also give further information, usually in a more elaborate syntactic form than the simpler pre-modifiers, such as a relative clause or a prepositional phrase (covered below). The noun phrases of the passage above can be parsed as shown in Table 2.16. The head nouns may be personal pronouns (*I, it*), demonstrative and relative pronouns (*those*), coordinated nouns (*sanction and consent*), or common nouns (*individual*). The modifiers are mostly adjectives (*impure*, *pure*) or verbal forms functioning as adjectives (*limited*). The post-modifiers are interesting, in that, unlike the (pre-)modifiers, they are typically full phrases themselves, rather than isolated words. They include relative clauses (which are a kind of dependent sentence, e.g., *[those] who know and can do better than I*), as well as prepositional phrases (*of the governed*).

**Table 2.16** *NP*'s of the sample passage.

| Np | Det | Mod | Head Noun | Post-Mod |
|----|-----|-----|-----------|----------|
| 1 | the | | **authority** | of government |
| 2 | | even | **such** | as I am willing to submit to |
| 3 | | | **I** | |
| 4 | | | **those** | who know and can do better than I |
| 5 | | many | **things** | |
| 6 | | even | **those** | who neither know nor can do so well |
| 7 | an | impure | **one** | |
| 8 | | | **it** | |
| 9 | the | | **sanction and consent** | of the governed |
| 10 | no | pure | **right** | over my person … concede to it. |
| 11 | the | | **progress** | from an absolute to a limited monarchy |
| 12 | an | absolute | **[monarchy]** | |
| 13 | a | limited | **monarchy** | |
| 14 | a | | **democracy** | |
| 15 | a | | **progress** | |
| 16 | a | true | **respect** | for the individual |
| 17 | the | | **individual** | |

Prepositions express spatial and temporal relations, among others. These are also said to project according to the *X*-template, but usually lack a pre-modifier. Some examples from the sample passage are listed in Table 2.17. The complements of PP are generally NP's, which may be simple head nouns like *government*. However, other complement types, such as the verb phrase in *after discussing it with Jo*, are also possible.

For verb phrases, the postmodifier (or complement) of a head verb would typically be one or more *NP* (noun phrase) maximal projections, which might, for example, function as a direct object in a *VP* like *pet the cat*. The complement may or may not be optional, depend-

ing on characteristics of the head. We can now make some language-specific generalizations about English. Some verbs, such as *give*, may take more than one kind of complement. So an appropriate template for a VP maximal projection in English would appear abstractly as *VP ⇒ (modifier) verb (modifier) (Complement1, Complement2 ComplementN)*. Complements are usually regarded as maximal projections, such as *NP*, *ADJP*, etc., and are enumerated in the template above, to cover possible multi-object verbs, such as *give*, which take both direct and indirect objects. Certain types of adverbs (really, quickly, smoothly, etc.) could be considered fillers for the VP modifier slots (before and after the head). In the sample passage, we find the following verb phrases as shown in Table 2.18.

**Table 2.17** PP's of the sample passage.

| Head Prep | Complement (Postmodifier) |
|---|---|
| **of** | Government |
| **as** | I am willing to submit to |
| **than** | I |
| **in** | many things |
| **of** | the governed |
| **over** | my person and property |
| **to** | it |
| **from** | an absolute [monarchy] |
| **to** | a limited monarchy |
| **to** | a democracy |
| **toward** | a true respect [for the individual] |
| **for** | the individual |

**Table 2.18** VP's of the sample passage.

| Pre-mod | Verb Head | Post-mod | Complement |
|---|---|---|---|
| | **submit to** | | [the authority of government] |
| cheerfully | **obey** | | those who know and can do better than I |
| | **is** | still | an impure one |
| | **be** | | strictly just |
| | **have** | | the sanction |
| | **have** | | no pure right |
| | **concede** | | to it |
| | **is** | | a progress |

    *VP* presents some interesting issues. First, notice the multi-word verb *submit to*. Multi-word verbs such as *look after* and *put up with* are common. We also observe a number of auxiliary elements clustering before the verb in sentences of the sample passage: *am willing to submit to*, *will cheerfully obey*, and *can do better*. Rather than considering these as simple modifiers of the verbal head, they can be taken to have *scope* over the *VP* as a whole, which implies they are outside the *VP*. Since they are outside the *VP*, we can assume them to be

heads in their own right, of phrases which require a *VP* as their complement. These elements mainly express tense (time or duration of verbal action) and modality (likelihood or probability of verbal action). In a full sentence, the *VP* has explicit or implicit inflection (projected from its verbal head) and indicates the person, number and other context-dependent features of the verb in relation to its arguments. In English, the person (first, second, third) and number (singular, plural) attributes, collectively called agreement features, of subject and verb must match. For simplicity, we will lump all these considerations together as inflectional elements, and posit yet another phrase type, the Inflectional Phrase (*IP*): *IP* $\Rightarrow$ *premodifier head VP-complement.*

The premodifier slot (sometimes called the *specifier* position in linguistic theory) of an *IP* is often filled by the subject of the sentence (typically a noun or *NP*). Since the *IP* unites the subject of a sentence with a *VP*, *IP* can also be considered simply as the sentence category, often written as *S* in speech grammars.

## 2.4.1.2. Clauses and Sentences

The *subject* of a sentence is what the sentence is mainly about. A *clause* is any phrase with both a subject and a VP (*predicate* in traditional grammars) that has potentially independent interpretation – thus, for us, a clause is an IP, a kind of sentence. A phrase is a constituent lacking either subject, predicate, or both. We have reviewed a number of phrase types above. There are also various types of clauses and sentences.

Even though clauses are sentences from an internal point of view (having subject and predicate), they often function as simpler phrases or words would, e.g., as modifiers (adjective and adverbs) or nouns and noun phrases. Clauses may appear as post-modifiers for nouns (so-called *relative clauses*), basically a kind of adjective clause, sharing their subjects with the containing sentence. Some clauses function as NP's in their own right. One common clause type substitutes a *wh-word* like *who* or *what* for a direct object of a verb in the embedded clause, to create a *questioned noun phrase* or indirect question: (*I don't know who Jo saw*.). In these clauses, it appears to syntacticians that the *questioned* object of the verb [VP saw **who**] has been extracted or moved to a new surface position (following the main clause verb *know*). This is sometimes shown in the phrase-structure diagram by co-indexing an empty *ghost* or trace constituent at the original position of the question pronoun with the question-NP appearing at the surface site:

*I don't know [$_{NPobj}$ [IP [$_{NPi}$ who] Jo saw [$_{NPi}$ _ ]]]*
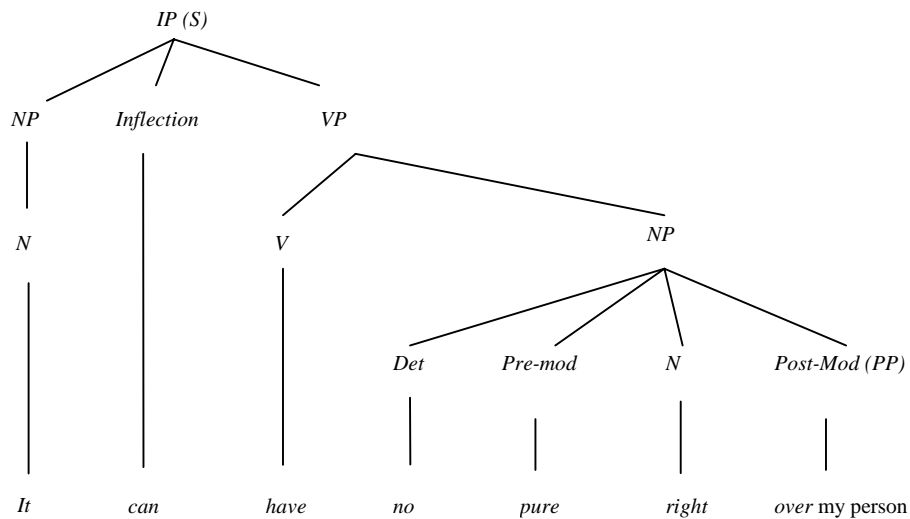*[$_{NPsubj}$ [$_{IP}$ Whoever wins the game]] is our hero.*

There are various characteristic types of sentences. Some typical types include:

Declarative: *I gave her a book.*

Yes-no question: *Did you give her a book* ?

Wh-question: *What did you give her*?

Alternatives question: *Did you give her a book, a scarf, or a knife*?

Tag question: *You gave it to her, didn't you?*

Passive: *She was given a book.*

Cleft: *It must have been a book that she got.*

Exclamative: *Hasn't this been a great birthday!*

Imperative: *Give me the book.*

## 2.4.1.3.    Parse Tree Representations

Sentences can be diagrammed in parse trees to indicate phrase-internal structure and linear precedence and immediate dominance among phrases. A typical phrase-structure tree for part of an embedded sentence is illustrated in Figure 2.26.



**Figure 2.26** A simplified phrase-structure diagram.

For brevity, the same information illustrated in the tree can be represented as a bracketed string as follows:

*[IP [ NP [ N It ]N]NP [I can ]I [VP [V have ]V [NP no pure right [PP over my person ]PP]NP]VP]IP*

With such a bracketed representation, almost every type of syntactic constituent can be coordinated or joined with another of its type, and usually a new phrase node of the common type is added to subsume the constituents such as *NP: We have [NP [NP tasty berries] and [NP tart juices]], IP/S: [IP [IP Many have come] and [IP most have remained]], PP: We went [PP [PP over the river] and [PP into the trees]],* and *VP: We want to [VP [VP climb the mountains] and [VP sail the seas]].*

## 2.4.2.    Semantic Roles

In traditional syntax, grammatical roles are used to describe the direction or control of action relative to the verb in a sentence. Examples include the ideas of *subject*, *object*, *indirect object*, etc. Semantic roles, sometimes called case relations, seem similar but dig deeper. They are used to make sense of the participants in an event, and they provide a vocabulary for us to answer the basic question *who did what to whom*. As developed by [13] and others, the theory of semantic roles posits a limited number of universal roles. Each basic meaning of each verb in our mental dictionary is tagged for the obligatory and optional semantic roles used to convey the particular meaning. A typical inventory of case roles is given below:

| | |
|---|---|
| *Agent* | *cause or initiator of action, often intentional* |
| *Patient/Theme* | *undergoer of the action* |
| *Instrument* | *how action is accomplished* |
| *Goal* | *to whom action is directed* |
| *Result* | *result of action* |
| *Location* | *location of action* |

These can be realized under various syntactic identities, and can be assigned to both required complement and optional adjuncts. A noun phrase in the Agentive role might be the surface subject of a sentence, or the object of the preposition *by* in a passive. For example, the verb *put* can be considered a process that has, in one of its senses, the case role specifications shown in Table 2.19.

**Table 2.19** Analysis of a sentence with *put*.

| Analysis | Example | | | |
|---|---|---|---|---|
| | *Kim* | *put* | *the book* | *on the table.* |
| ***Grammatical functions*** | *Subject (NP)* | *Predicate (VP)* | *Object (NP)* | *Adverbial (ADVP)* |
| ***Semantic roles*** | *Agent* | *Instrument* | *Theme* | *Location* |

Now consider this passive-tense example, where the semantic roles align with different grammatical roles shown in Table 2.20. Words that look and sound identical can have different meaning or different *senses* as shown in

Table **2.21**. The sporting sense of *put* (as in the sport of shot-put), illustrates the meaning/sense-dependent nature of the role patterns, because in this sense the Locative case is no longer obligatory, as it is in the original sense illustrated in Table 2.19 and Table 2.20.

**Table 2.20** Analysis of passive sentence with *put*.

| Analysis | Example | | |
|---|---|---|---|
| | *The book* | *was put* | *on the table* |
| ***Grammatical functions*** | *Subject (NP)* | *Predicate (VP)* | *Adverbial (ADVP)* |
| ***Semantic roles*** | *Agent* | *Instrument* | *Location* |

**Table 2.21** Analysis of a different pattern of *put.*

| Analysis | Example | | |
|---|---|---|---|
| | *Kim* | *put* | *the shot.* |
| *Grammatical functions* | *Subject (NP)* | *Predicate (VP)* | *Object (NP)* |
| *Semantic Roles* | *Agent* | *Instrument* | *Theme* |

The lexical meaning of a verb can be further decomposed into primitive semantic relations such as CAUSE, CHANGE, and BE. The verb *open* might appear as *CAUSE(NP1,PHYSICAL-CHANGE(NP2,NOT-OPEN,OPEN)).* This says that for an agent (NP1) to *open* a theme (NP2) is to cause the patient to change from a not-opened state to an opened state. Such systems can be arbitrarily detailed and exhaustive, as the application requires.

## 2.4.3.      Lexical Semantics

The specification of particular meaning templates for individual senses of particular words is called *lexical semantics*. When words combine, they may take on propositional meanings resulting from the composition of their meanings in isolation. We could imagine that a speaker starts with a proposition in mind (logical form as will be discussed in the next section), creating a need for particular words to express the idea (lexical semantics); the proposition is then linearized (syntactic form) and spoken (phonological/phonetic form). Lexical semantics is the level of meaning before words are composed into phrases and sentences, and it may heavily influence the possibilities for combination.

Words can be defined in a large number of ways including by relations to other words, in terms of decomposition semantic primitives, and in terms of non-linguistic cognitive constructs, such as perception, action, and emotion. There are hierarchical and non-hierarchical relations. The main hierarchical relations would be familiar to most object-oriented programmers. One is *is-a* taxonomies (a *crow* is-a *bird*), which have transitivity of properties from type to subtype (inheritance). Another is *has-a* relations (a *car* has-a *windshield*), which are of several differing qualities, including process/subprocess (*teaching* has-a subprocess *giving exams*), and arbitrary or natural subdivisions of part-whole relations (*bread* has-a division into *slices, meter* has-a division into *centimeters*). Then there are non-branching hierarchies (no fancy name) that essentially form scales of degree, such as *frozen* ⇒ *cold* ⇒ *lukewarm* ⇒ *hot* ⇒ *burning*. Non-hierarchical relations include synonyms, such as *big/large*, and antonyms such as *good/bad*.

Words seem to have natural affinities and disaffinities in the semantic relations among the concepts they express. Because these affinities could potentially be exploited by future language understanding systems, researchers have used the generalizations above in an attempt to tease out a parsimonious and specific set of basic relations under which to group entire lexicons of words. A comprehensive listing of the families and subtypes of possible semantic relations has been presented in [10]. In Table 2.22, the leftmost column shows names for families of proposed relations, the middle column differentiates subtypes within

each family, and the rightmost column provides examples of word pairs that participate in the proposed relation. Note that case roles have been modified for inclusion as a type of semantic relation within the lexicon.

**Table 2.22** Semantic relations.

| Family | Subtype | Example |
|---|---|---|
| Contrasts | Contrary | *old-young* |
| | Contradictory | *alive-dead* |
| | Reverse | *buy-sell* |
| | Directional | *front-back* |
| | Incompatible | *happy-morbid* |
| | Asymmetric contrary | *hot-cool* |
| | Attribute similar | *rake-fork* |
| Similars | Synonymity | *car-auto* |
| | Dimensional similar | *smile-laugh* |
| | Necessary attribute | *bachelor-unmarried* |
| | Invited attribute | *food-tasty* |
| | Action subordinate | *talk-lecture* |
| Class Inclusion | Perceptual subord. | *animal-horse* |
| | Functional subord. | *furniture-chair* |
| | State subord. | *disease-polio* |
| | Activity subord. | *game-chess* |
| | Geographic subord. | *country-Russia* |
| | Place | *Germany-Hamburg* |
| Case Relations | Agent-action | *artist-paint* |
| | Agent-instrument | *farmer-tractor* |
| | Agent-object | *baker-bread* |
| | Action-recipient | *sit-chair* |
| | Action-instrument | *cut-knife* |
| Part-Whole | Functional object | *engine-car* |
| | Collection | *forest-tree* |
| | Group | *choir-singer* |
| | Ingredient | *table-wood* |
| | Functional location | *kitchen-stove* |
| | Organization | *college-admissions* |
| | Measure | *mile-yard* |

We can see from Table 2.22 that a single word could participate in multiple relations of different kinds. For example, *knife* appears in the examples for *Similars: invited attribute* (i.e., a desired and expected property) as: *knife-sharp*, and also under *Case Relations*: *ac-*

*tion-instrument*, which would label the relation of *knife* to the action *cut* in *He cut the bread with a knife.* This suggests that an entire lexicon could be viewed as a graph of semantic relations, with words or idioms as nodes and connecting edges between them representing semantic relations as listed above. There is a rich tradition of research in this vein.

The biggest practical problem of lexical semantics is the context-dependent resolution of senses of words – so-called polysemy. A classic example is *bank - bank of the stream* as opposed to *money in the bank*. While lexicographers try to identify distinct senses when they write dictionary entries, it has been generally difficult to rigorously quantify exactly what counts as a discrete sense of a word and to disambiguate the senses in practical contexts. Therefore, designers of practical speech understanding systems generally avoid the problem by limiting the domain of discourse. For example, in a financial application, generally only the sense of *bank* as a fiduciary institution is accessible, and others are assumed not to exist. It is sometimes difficult to make a principled argument as to how many distinct senses a word has, because at some level of depth and abstraction, what might appears as separate senses seem to be similar or related, as *face* could be *face of a clock* or *face of person*.

Senses are usually distinguished within a given part-of-speech (POS) category. Thus, when an occurrence of *bank* has been identified as a verb, the *shore* sense might be automatically eliminated, though depending on the sophistication of the system's lexicon and goals, there can be sense differences for many English verbs as well. Within a POS category, often the words that occur near a given ambiguous form in the utterance or discourse are clues to interpretation, where links can be established using semantic relations as described above. Mutual information measures as discussed in Chapter 3 can sometimes provide hints. In a context of dialog where other, less ambiguous financial terms come up frequently, the sense of *bank* as fiduciary institution is more likely. Finally, when all else fails, often senses can be ranked in terms of their a priori likelihood of occurrence. It should always be borne in mind that language is not static; it can change form under a given analysis at any time. For example, the stable English form *spinster*, a somewhat pejorative term for an older, never-married female, has recently taken on a new morphologically complex form, with the new sense of a high political official, or media spokesperson, employed to provide bland disinformation (*spin*) on a given topic.

### 2.4.4. Logical Form

Because of all the lexical, syntactic, and semantic ambiguity in language, some of which requires external context for resolution, it is desirable to have a metalanguage in which to concretely and succinctly express all linguistically possible meanings of an utterance before discourse and world knowledge are applied to choose the most likely interpretation. The favored metalanguage for this purpose is called the predicate logic, used to represent the logical form, or context-independent meaning, of an utterance. The semantic component of many SLU architectures builds on a substrate of two-valued, first-order, logic. To distinguish *shades of meaning* beyond truth and falsity requires more powerful formalisms for knowledge representation.

In a typical first-order system, predicates correspond to events or conditions denoted by verbs (such as *Believe* or *Like*), states of identity (such as being a *Dog* or *Cat*), and prop-

erties of varying degrees of permanence (*Happy*). In this form of logical notation, predicates have open places, filled by arguments, as in a programming language subroutine definition. Since individuals may have identical names, subscripting can be used to preserve unique reference. In the simplest systems, predication ranges over individuals rather than higher-order entities such as properties and relations.

Predicates with filled argument slots map onto sets of individuals (constants) in the universe of discourse, in particular those individuals possessing the properties, or participating in the relation, named by the predicate. One-place predicates like *Soldier*, *Happy*, or *Sleeps* range over sets of individuals from the universe of discourse. Two-place predicates, like transitive verbs such as *loves*, range over a set consisting of ordered pairs of individual members (constants) of the universe of discourse. For example, we can consider the universe of discourse to be $U$ = {*Romeo, Juliet, Paris, Rosaline, Tybalt}*, people as characters in a play. They do things with and to one another, such as loving and killing. Then we could imagine the relation *Loves* interpreted as the set of ordered pairs: {*<Romeo, Juliet>, <Juliet, Romeo>, <Tybalt, Tybalt>, <Paris, Juliet>*}, a subset of the Cartesian product of theoretically possible love matches $U \times U$ . So, for any ordered pair *x, y* in *U*, *Loves*(*x, y*) is true iff the ordered pair *<x,y>* is a member of the extension of the *Loves* predicate as defined, e.g., *Romeo loves Juliet*, *Juliet loves Romeo*, etc.. Typical formal properties of relations are sometimes specially marked by grammar, such as the reflexive relation *Loves*(*Tybalt*, *Tybalt*), which can rendered in natural language as *Tybalt loves himself*. Not every possibility is present; for instance in our example, the individual Rosaline does not happen to participate at all in this extensional definition of *Loves* over *U*, as her omission from the pairs list indicates. Notice that the subset of *Loves*(*x, y*) of ordered pairs involving both *Romeo* and *Juliet* is symmetric, also marked by grammar, as in *Romeo and Juliet love each other.* This general approach extends to predicates with any arbitrary number of arguments, such as intransitive verbs like *give*.

Just as in ordinary propositional logic, connectives such as negation, conjunction, disjunction, and entailment are admitted, and can be used with predicates to denote common natural language meanings:

> Romeo isn't happy = ^Happy(Romeo)
> Romeo isn't happy, but Tybalt is (happy) = ^Happy(Romeo) && Happy(Tybalt)
> Either Romeo or Tybalt is happy = Happy(Romeo) || Happy(Tybalt)
> If Romeo is happy, Juliet is happy = Happy(Romeo) ➔ Happy(Juliet)

Formulae, such as those above, are also said to bear a binary truth value, true or false, with respect to a world of individuals and relations. The determination of the truth value is compositional, in the sense that the truth value of the whole depends on the truth value of the parts. This is a simplistic but formally tractable view of the relation between language and meaning.

Predicate logic can also be used to denote quantified noun phrases. Consider a simple case such as *Someone killed Tybalt*, predicated over our same $U$ = {*Romeo, Juliet, Paris, Rosaline, Tybalt*}. We can now add an *existential* quantifier, , standing for *there exists* or *there is at least one*. This quantifier will bind a variable over individuals in *U*, and will attach to a proposition to create a new, quantified proposition in logical form. The use of vari-

ables in propositions such as *killed(x, y)* creates open propositions. Binding the variables with a quantifier over them closes the proposition. The quantifier is prefixed to the original proposition:

   *x Killed(x, Tybalt)*

To establish a truth (semantic) value for the quantified proposition, we have to satisfy the disjunction of propositions in *U*: *Killed*(*Romeo*, *Tybalt*) ∨ *Killed*(*Juliet, Tybalt*) ∨ *Killed*(*Paris, Tybalt*) ∨ *Killed*(*Rosaline, Tybalt*) ∨ *Killed*(*Tybalt, Tybalt*). The set of all such bindings of the variable *x* is the space that determines the truth or falsity of the proposition. In this case, the binding of *x = Romeo* is sufficient to assign a value true to the existential proposition.

## 2.5.    HISTORICAL PERSPECTIVE AND FURTHER READING

Motivated to improve speech quality over the telephone, AT&T Bell Labs has contributed many influential discoveries in speech hearing, including the critical band and articulation index [2, 3]. The *Auditory Demonstration* CD prepared by Houtsma, Rossing, and Wagenaars [18] has a number of very interesting examples on psychoacoustics and its explanations. *Speech, Language, and Communication* [30] and *Speech Communication - Human and Machines* [32] are two good books that provide modern introductions to the structure of spoken language. Many speech perception experiments were conducted by exploring how phonetic information is distributed in the time or frequency domain. In addition to the formant structures for vowels, frequency importance function [12] has been developed to study how features related to phonetic categories are stored at various frequencies. In the time domain, it has been observed [16, 19, 42] that salient perceptual cues may not be evenly distributed over the speech segments and that certain perceptual critical points exist.

As intimate as speech and acoustic perception may be, there are also strong evidences that lexical and linguistic effects on speech perception are not always consistent with acoustic ones. For instance, it has long been observed that humans exhibit difficulties in distinguishing non-native phonemes. Human subjects also carry out categorical goodness difference assimilation based on their mother tongue [34], and such perceptual mechanism can be observed as early as in six-month-old infants [22]. On the other hand, hearing-impaired listeners are able to effortlessly overcome their acoustical disabilities for speech perception [8]. Speech perception is not simply an auditory matter. McGurk and MacDonald (1976) [27, 28] dramatically demonstrated this when they created a video tape on which the auditory information (phonemes) did not match the visual speech information. The effect of this mismatch between the auditory signal and the visual signal was to create a third phoneme different from both the original auditory and visual speech signals. An example is dubbing the phoneme /*ba*/ to the visual speech movements /*ga*/. This mismatch results in hearing the phoneme /*da*/. Even when subjects know of the effect, they report the McGurk effect percept. The McGurk effect has been demonstrated for consonants, vowels, words, and sentences.

The earliest scientific work on phonology and grammars goes back to Panini, a Sanskrit grammarian of the fifth century B.C. (estimated), who created a comprehensive and

scientific theory of phonetics, phonology, and morphology, based on data from Sanskrit (the classical literary language of the ancient Hindus). Panini created formal production rules and definitions to describe Sanskrit grammar, including phenomena such as construction of sentences, compound nouns, etc. Panini's formalisms function as ordered rules operating on underlying structures in a manner analogous to modern linguistic theory. Panini's phonological rules are equivalent in formal power to Backus-Nauer form (BNF). A general introduction to this pioneering scientist is Cardona [9].

An excellent introduction to all aspects of phonetics is *A Course in Phonetics* [24]. A good treatment of the acoustic structure of English speech sounds and a through introduction and comparison of theories of speech perception is to be found in [33]. The basics of phonology as part of linguistic theory are treated in *Understanding Phonology* [17]. An interesting treatment of word structure (morphology) from a computational point of view can be found in *Morphology and Computation* [40]. A comprehensive yet readable treatment of English syntax and grammar can be found in *English Syntax* [4] and *A Comprehensive Grammar of the English Language* [36]. Syntactic theory has traditionally been the heart of linguistics, and has been an exciting and controversial area of research since the 1950s. Be aware that almost any work in this area will adopt and promote a particular viewpoint, often to the exclusion or minimization of others. A reasonable place to begin with syntactic theory is *Syntax: A Minimalist Introduction* [37]. An introductory textbook on syntactic and semantic theory that smoothly introduces computational issues is *Syntactic Theory: A Formal Introduction* [39]. For a philosophical and entertaining overview of various aspects of linguistic theory, see *Rhyme and Reason: An Introduction to Minimalist Syntax* [44]. A good and fairly concise treatment of basic semantics is *Introduction to Natural Language Semantics* [11]. Deeper issues are covered in greater detail and at a more advanced level in *The Handbook of Contemporary Semantic Theory* [25]). The intriguing area of lexical semantics (theory of word meanings) is comprehensively presented in *The Generative Lexicon* [35]. *Concise History of the Language Sciences* [21] is a good edited book if you are interested in the history of linguistics.

## REFERENCES

[1]     Aliprand, J., *et al.*, *The Unicode Standard, Version 2.0*, 1996, Addison Wesley.

[2]     Allen, J.B., "How Do Humans Process and Recognize Speech?," *IEEE Trans. on Speech and Audio Processing*, 1994, **2**(4), pp. 567-577.

[3]     Allen, J.B., "Harvey Fletcher 1884--1981" in *The ASA Edition of Speech and Hearing Communication* 1995, Woodbury, New York, pp. A1-A34, Acoustical Society of America.

[4]     Baker, C.L., *English Syntax*, 1995, Cambridge, MA, MIT Press.

[5]     Blauert, J., *Spatial Hearing*, 1983, MIT Press.

[6]     Brill, E., "Transformation-Based Error-Driven Learning and Natural Language Processing: A Case Study in Part-of-Speech Tagging," *Computational Linguistics*, 1995, **21**(4), pp. 543-566.

[7]     Brown, P., *et al.*, "Class-Based N-gram Models of Natural Language," *Computational Linguistics*, 1992, **18**(4).

[8]     Caplan, D. and J. Utman, "Selective Acoustic Phonetic Impairment and Lexical Access in an Aphasic Patient," *Journal of the Acoustical Society of America*, 1994, **95**(1), pp. 512-517.

[9]     Cardona, G., *Panini: His Work and Its Traditions: Background and Introduction*, 1988, Motilal Banarsidass.

[10]    Chaffin, R., Herrmann, D., "The Nature of Semantic Relations: A Comparison of Two Approaches" in *Representing knowledge in Semantic Networks*, M. Evens, Editor 1988, Cambridge, UK, Cambridge University Press.

[11]    de Swart, H., *Introduction to Natural Language Semantics*, 1998, Stanford, California, USA, Center for the Study of Language and Information Publications.

[12]    Duggirala, V., *et al.*, "Frequency Importance Function for a Feature Recognition Test Material," *Journal of the Acoustical Society of America*, 1988, **83**(9), pp. 2372-2382.

[13]    Fillmore, C.J., "The Case for Case" in *Universals in Linguistic Theory*, E. Bach and R. Harms, eds. 1968, New York, NY, Holt, Rinehart and Winston.

[14]    Fletcher, H., "Auditory patterns," *Rev. Mod. Phys.*, 1940, **12**, pp. 47-65.

[15]    Fry, D.B., *The Physics of Speech*, Cambridge Textbooks in Linguistics, 1979, Cambridge, U.K., Cambridge University Press.

[16]    Furui, S., "On The Role of Spectral Transition for Speech Perception," *Journal of the Acoustical Society of America*, 1986, **80**(4), pp. 1016-1025.

[17]    Gussenhoven, C., Jacobs, H., *Understanding Phonology*, Understanding Language Series, 1998, Edward Arnold.

[18]    Houtsma, A., T. Rossing, and W. Wagenaars, *Auditory Demonstrations*, 1987, Institute for Perception Research, Eindhovern, The Netherlands, Acoustic Society of America.

[19]    Jenkins, J., W. Strange, and S. Miranda, "Vowel Identification in Mixed-Speaker Silent-Center Syllables," *Journal of the Acoustical Society of America*, 1994, **95**(2), pp. 1030-1041.

[20]    Klatt, D., "Review of the ARPA Speech Understanding Project," *Journal of Acoustical Society of America*, 1977, **62**(6), pp. 1324-1366.

[21]    Koerner, E. and E. Asher, eds. *Concise History of the Language Sciences*, , 1995, Oxford, Elsevier Science.

[22]    Kuhl, P., "Infant's Perception and Representation of Speech: Development of a New Theory," *Int. Conf. on Spoken Language Processing*, 1992, Alberta, Canada pp. 449-452.

[23]    Kupeic, J., "Robust Part-of-Speech Tagging Using a Hidden Markov Model," *Computer Speech and Language*, 1992, **6**, pp. 225-242.

[24]    Ladefoged, P., *A Course in Phonetics*, 1993, Harcourt Brace Johanovich.

[25]    Lappin, S., *The Handbook of Contemporary Semantic Theory*, Blackwell Handbooks in Linguistics, 1997, Oxford, UK, Blackwell Publishsers Inc.

[26]    Lindsey, P. and D. Norman, *Human Information Processing*, 1972, New York and London, Academic Press.

[27]    MacDonald, J. and H. McGurk, "Visual Influence on Speech Perception Process," *Perception and Psychophysics*, 1978, **24**(3), pp. 253-257.

[28]    McGurk, H. and J. MacDonald, "Hearing Lips and Seeing Voices," *Nature*, 1976, **264**, pp. 746-748.

[29]    Merialdo, B., "Tagging English Text with a Probabilistic Model," *Computational Linguistics*, 1994, **20**(2), pp. 155-172.

[30]    Miller, J. and P. Eimas, *Speech, Language and Communication*, Handbook of Perception and Cognition, eds. E. Carterette and M. Friedman, 1995, Academic Press.

[31]    Moore, B.C., *An Introduction to the Psychology of Hearing*, 1982, London, Academic Press.

[32]    O'Shaughnessy, D., *Speech Communication -- Human and Machine*, 1987, Addison-Wesley.

[33]    Pickett, J.M., *The Acoustics of Speech Communication*, 1999, Needham Heights, MA, Allyn & Bacon.

[34]    Polka, L., "Linguistic Influences in Adult Perception of Non-native Vowel Contrast," *Journal of the Acoustical Society of America*, 1995, **97**(2), pp. 1286-1296.

[35]    Pustejovsky, J., *The Generative Lexicon*, 1998, Bradford Books.

[36]    Quirk, R., Svartvik, J., Leech, G., *A Comprehensive Grammar of the English Language*, 1985, Addison-Wesley Pub Co.

[37]    Radford, A., *Syntax: A Minimalist Introduction*, 1997, Cambridge, U.K., Cambridge Univ. Press.

[38]    Rossing, T.D., *The Science of Sound*, 1982, Reading, MA, Addison-Wesley.

[39]    Sag, I., Wasow, T., *Syntactic Theory: A Formal Introduction*, 1999, Cambridge, UK, Cambridge University Press.

[40]    Sproat, R., *Morphology and Computation*, ACL-MIT Press Series in Natural Language Processing, 1992, Cambridge, MA, MIT Press.

[41]    Stevens, S.S. and J. Volkman, "The Relation of Pitch to Frequency," *Journal of Psychology*, 1940, **53**, pp. 329.

[42]    Strange, W., J. Jenkins, and T. Johnson, "Dynamic Specification of Coarticulated Vowels," *Journal of the Acoustical Society of America*, 1983, **74**(3), pp. 695-705.

[43]    Thoreau, H.D., *Civil Disobedience, Solitude and Life Without Principle*, 1998, Prometheus Books.

[44]    Uriagereka, J., *Rhyme and Reason: An Introduction to Minimalist Syntax*, 1998, Cambridge, MA, MIT Press.

[45]    Voutilainen, A., "Morphological Disambiguation" in *Constraint Grammar: A Language-Independent System for Parsing Unrestricted Text* 1995, Berlin, Mouton de Gruyter.

[46]    Weischedel, R., "BBN: Description of the PLUM System as Used for MUC-6," *The 6th Message Understanding Conferences (MUC-6)*, 1995, San Francisco, Morgan Kaufmann pp. 55-70.